

**Evidence of pervasive biologically functional secondary-structures within the  
genomes of eukaryotic single-stranded DNA viruses**

Brejnev Muhizi Muhire<sup>1</sup>, Michael Golden<sup>1</sup>, Ben Murrell<sup>2</sup>, Pierre Lefeuve<sup>3</sup>, Jean-Michel  
Lett<sup>3</sup>, Alistair Gray<sup>4</sup>, Art YF Poon<sup>5,6</sup>, Nobubelo Kwanele Ngandu<sup>7</sup>, Yves Semegni<sup>8</sup>, Emil  
Pavlov Tanov<sup>9</sup>, Adérito Luis Monjane<sup>1,4</sup>, Gordon William Harkins<sup>9</sup>, Arvind Varsani<sup>10,11,12</sup>,  
Dionne Natalie Shepherd<sup>4</sup>, Darren Patrick Martin<sup>1,#</sup>

<sup>1</sup>Institute of Infectious Diseases and Molecular Medicine, Computational Biology Group,  
University of Cape Town, Cape Town, South Africa

<sup>2</sup>Department of Medicine, University of California, San Diego, California, USA

<sup>3</sup>CIRAD, UMR PVBMT CIRAD-Université de la Réunion, Pôle de Protection des Plantes,  
Saint-Pierre, La Réunion, France

<sup>4</sup>Department of Molecular and Cell Biology, University of Cape Town, Rondebosch, Cape  
Town, South Africa

<sup>5</sup>BC Centre for Excellence in HIV/AIDS, Vancouver, Canada

<sup>6</sup>Department of Medicine, University of British Columbia, Vancouver, Canada

<sup>7</sup>Institute of Infectious Diseases and Molecular Medicine, Division of Medical Virology,  
University of Cape Town, Cape Town, South Africa

<sup>8</sup>Department of Mathematics and Physics, Cape Peninsula University of Technology, Cape  
Town, South Africa

<sup>9</sup>South African National Bioinformatics Institute, University of the Western Cape, Cape  
Town, South Africa

24 <sup>10</sup>School of Biological Sciences, University of Canterbury, Christchurch, New Zealand  
 25 <sup>11</sup>Biomolecular Interaction Centre, University of Canterbury, Christchurch, New Zealand  
 26 <sup>12</sup>Electron Microscope Unit, Division of Medical Biochemistry, Department of Clinical  
 27 Laboratory Sciences, University of Cape Town, Rondebosch, Cape Town, South Africa  
 28  
 29 Current address: Institute of Infectious Diseases and Molecular Medicine, Computational  
 30 Biology Group, University of Cape Town, Cape Town, South Africa  
 31 Corresponding Author: [darrenpatrickmartin@gmail.com](mailto:darrenpatrickmartin@gmail.com)  
 32 Running title: Pervasive secondary-structures in ssDNA virus genomes  
 33 Abstract word count: 242  
 34 Text word count: 7520  
 35

36 **Abstract**

37 Single-stranded DNA (ssDNA) viruses have genomes that are potentially capable of  
38 forming complex secondary-structures through Watson-Crick base-pairing between their  
39 constituent nucleotides. A few of the structural elements formed by such base-pairings  
40 are, in fact, known to have important functions during the replication of many ssDNA  
41 viruses. What is unknown, however, is (i) whether numerous additional ssDNA virus  
42 genomic structural elements predicted to exist by computational DNA folding methods  
43 actually exist, and (ii) whether those structures that do exist have any biological relevance.  
44 We therefore computationally inferred lists of the most evolutionarily conserved structures  
45 within a diverse selection of animal- and plant-infecting ssDNA viruses drawn from the  
46 families *Circoviridae*, *Anelloviridae*, *Parvoviridae*, *Nanoviridae* and *Geminiviridae*, and  
47 analysed these for evidence of natural selection favouring the maintenance of these  
48 structures. While we find evidence that is consistent with purifying selection being stronger  
49 at nucleotide sites that are predicted to be base-paired than it is at sites predicted to be  
50 unpaired, we also find strong associations between sites that are predicted to pair with one  
51 another and site pairs that are apparently coevolving in a complementary fashion.  
52 Collectively, these results indicate that natural selection actively preserves much of the  
53 pervasive secondary-structure that is evident within eukaryote-infecting ssDNA virus  
54 genomes and, therefore, that much of this structure is biologically functional. Lastly, we  
55 provide examples of various highly conserved but completely uncharacterised structural  
56 elements that likely have important functions within some of the ssDNA virus genomes  
57 analysed here.

58

59 **Introduction**

60 Besides encoding structural, regulatory and enzymatic proteins, the nucleotide sequences  
61 of viral genomes encode a wide range of regulatory motifs associated with, amongst other  
62 things, transcription (1, 2), translation (3), replication (4) and genome packaging (5). Other  
63 types of biologically relevant information encoded by many nucleotide sequences,  
64 including those of viruses, are the thermodynamically stable secondary and tertiary  
65 structures that these sequences form under physiological conditions.

66

67 While the capacity of single-stranded nucleic acid molecules to fold into higher order  
68 structures is crucial in all living organisms for the correct functioning of transfer RNA,  
69 ribosomal RNA, messenger RNA and small regulatory RNA molecules, such structures  
70 are also particularly important in the biology of many viruses with single-stranded DNA  
71 (ssDNA) and single-stranded RNA (ssRNA) genomes. Such structures can play vital roles  
72 during the entire viral reproductive cycle including the initiation of genome replication (6–  
73 12), the regulation of gene expression (13), the control of transcription (14), translation  
74 (15–17) and gene splicing (18), and the modulation of host anti-viral responses (19–21).

75

76 Within viral genomes, biologically important structural elements tend to be highly  
77 conserved across even distantly related species (22). In ssRNA viruses, for example, they  
78 include the rev response element (RRE) in Human and Simian immunodeficiency viruses  
79 (23–25), the *cis*-acting replication elements (CREs) of flaviviruses (26), luteoviruses (27),  
80 carmoviruses (11), coronaviruses (28), alphaflexiviruses (29), reoviruses (30), and  
81 picornaviruses (31, 32), the internal ribosomal entry site (IRES) elements of flaviviruses

82 (33, 34), picornaviruses (35), pestiviruses (36) and dicistroviruses (37), and the cap-  
83 independent translation elements (CITEs) found in many plant-infecting ssRNA viruses  
84 (38, 39).

85

86 Similarly, in many ssDNA virus genomes DNA secondary-structural elements have been  
87 identified that have crucial biological functions. While in parvoviruses these include  
88 structural elements that are essential for genome replication (9, 40, 41) and optimal gene  
89 expression (42–46), in geminiviruses, nanoviruses and circoviruses highly conserved  
90 stem-loop structures at their replication origins are essential for the initiation and  
91 termination of replication (6, 47–49). Besides these few examples, however, it is currently  
92 unknown how pervasive biologically important secondary-structural elements are within  
93 the genomes of these viruses (50–52).

94

95 It is important to stress the distinction between the simple existence within viral genomes  
96 of pervasive stable secondary-structures and the biological importance of these structures.  
97 Above a certain length even randomly generated single-stranded oligonucleotides will form  
98 stable secondary-structures (53), and it is therefore plausible that many essentially  
99 functionless secondary-structural elements might exist within ssDNA and ssRNA viral  
100 genomes.

101

102 It is, however, theoretically possible to computationally determine the functional  
103 importance within viral genomes of secondary-structural elements (detected either by  
104 computational prediction or by more rigorous laboratory analyses) by simply examining

105 patterns of evolution that are evident within groups of related genome sequences.  
106 Specifically, although biologically functional secondary-structural elements should be  
107 evolutionarily conserved across diverse viral lineages, the nucleotide sequences from  
108 which these elements are composed should display distinctive signals of natural selection  
109 favouring the maintenance of these structures. Whereas in coding regions these signals  
110 might include codon usage biases (54, 55), and decreased rates of synonymous  
111 substitution (56), throughout the genome these signals could also include high rates of  
112 reversion substitution (51, 57) and increased frequencies of complementarily coevolving  
113 nucleotide pairs – particularly amongst those nucleotides predicted to be base-paired  
114 within secondary-structural elements (58–60).

115

116 Accordingly, experimental investigations of individual structural elements within some  
117 ssDNA virus genomes have clearly demonstrated the existence of strong natural selection  
118 favouring the maintenance of these elements. For example, when mutations were  
119 experimentally introduced that disrupted particular base-pairings within a stem-loop  
120 structure at the origin of replication of the circovirus, *Porcine circovirus* 1 (PCV-1), the  
121 disrupted base-pairings were rapidly restored during replication through a DNA  
122 polymerase mediated template switching mechanism (61, 62). Similarly, in the  
123 geminivirus, *Maize streak virus* (MSV), mutations that potentially disrupted base-pairings  
124 within a complex computationally predicted structural element were found to very  
125 predictably revert to the original nucleotide so as to re-stabilise the structural element (51).

126

127 Here we examine the biological relevance of pervasive computationally predicted  
 128 secondary-structures within diverse eukaryote-infecting ssDNA virus genomes. After using  
 129 a free energy minimisation approach to identify conserved secondary-structural elements  
 130 within groups of closely related full genome sequences, we applied various tests to  
 131 determine whether mutational processes differed between structured and unstructured  
 132 genome regions in ways consistent with the evolutionary conservation of the identified  
 133 structural elements. While we provide strong evidence of extensive biologically relevant  
 134 secondary-structures within eukaryotic ssDNA virus genomes, we further identify what are  
 135 likely some of the most functionally important uncharacterised structural elements within  
 136 these genomes.

137

## 138 **Materials and methods**

### 139 **Dataset preparation**

140 All available circovirus, anellovirus, parvovirus, nanovirus and geminivirus full genome  
 141 sequences were obtained from GenBank (<http://www.ncbi.nlm.nih.gov/pubmed/>) between  
 142 April 2011 and August 2012. Full genome sequences for each of the five families were  
 143 preliminarily aligned separately using MUSCLE (63) implemented in MEGA5 (64) and  
 144 subdivided into datasets of sequences sharing at least 75% sequence identity. This was  
 145 done to ensure reasonable alignment accuracy during subsequent sequence analyses  
 146 (65) while at the same time, providing enough sequence diversity to enable the accurate  
 147 characterisation of evolutionary processes acting to maintain predicted secondary-  
 148 structures.

149

150 A set of 23 datasets was obtained, each containing between 21 and 519 full genome  
151 sequences. Each of these full genome sequence datasets was realigned using MUSCLE  
152 (with default settings) and, where necessary, manually edited. The resulting alignments  
153 will hereafter be referred to as “large datasets” (Table 1). This distinction is important  
154 because many of the analyses performed could only be carried out on subsets of these  
155 datasets. Specifically, from each of the large datasets we first extracted an “intermediate  
156 dataset”. In all but four cases each contained one representative sequence from each of  
157 the 30 most divergent viral sequence lineages in the large datasets. The exceptional  
158 cases were the AnelloTTSuV1, AnelloTTV, ParvoHBoV, and ParvoMPV datasets that  
159 respectively contained only 21, 22, 21 and 26 sequences, all of which were included in the  
160 intermediate dataset. From each of the intermediate datasets we further extracted a “small  
161 dataset” containing one representative sequence of each of the 10 most divergent  
162 lineages.

163

#### 164 **Detection of conserved secondary-structural elements within ssDNA virus genomes**

165 Since biologically relevant secondary-structural elements are likely to be at least partially  
166 conserved during evolution, we used the computer program Nucleic Acid Secondary  
167 Structure Predictor (NASP; (66)), to identify the conserved secondary-structural elements  
168 within the set of representative full genomic sequences in each of the small datasets.

169

170 NASP takes as input a set of aligned nucleic acid sequences and uses Gibbs free-energy  
171 (67) and Boltzmann probability (68) techniques implemented in the hybrid-ss component  
172 of the UNAFOLD software package (69) to determine an ensemble of nearly minimum free



173 energy (MFE) folds for each of the input sequences. It then uses a nucleotide-shuffling  
174 based permutation test to statistically determine the sets of conserved structural elements  
175 within the folded sequences that contribute most to their over-all thermodynamic stability.

176

177 Precisely, NASP produces sets of “pairing matrices” for each of the input sequences in  
178 each small dataset which are then compressed into a “consensus base-pairing matrix,”  
179 called the M matrix, using a weighted sum of the pairing matrices obtained for each of the  
180 input sequences (66). The use by NASP of weighted sums in the calculation of its M  
181 Matrix is intended to counteract unavoidable sampling biases in sequence datasets so as  
182 to ensure that similar structures within very closely related sequences do not make unfair  
183 contributions to the “conservation scores” that NASP calculates for the individual structural  
184 elements that it identifies.

185

186 Importantly, in our study NASP provided a conservation score for each discrete structural  
187 element identified within the M-matrices calculated for each of the small datasets, and  
188 indicated the subsets of structural elements referred to as “high confidence structure sets”  
189 (HCSSs), that accounted for the analysed sequences having significantly lower MFE  
190 scores than those expected in randomised sequences with identical base compositions  
191 (66). Whereas the conservation scores for the individual structural elements provided an  
192 obvious way of ranking these in order of their likely biological relevance, the demarcation  
193 of HCSSs provided an objective means of focusing further analyses into the biological  
194 relevance of secondary-structures on just the structural elements that are most likely to  
195 really occur.

196

197 In our NASP analysis, sequences were folded as either linear (for the three parvovirus  
198 small datasets) or circular (for all 20 other small datasets) ssDNA at either 37°C (for  
199 animal-infecting circoviruses, anelloviruses and parvoviruses) or 25°C (for plant-infecting  
200 geminiviruses and nanoviruses) under 0M magnesium and 1M sodium ionic conditions.  
201 The HCSS was identified using 100 nucleotide shuffling permutations with a permutation  
202 *p*-value threshold of 0.05. In all subsequent analyses, the only nucleotides considered as  
203 being paired within secondary-structures were those occurring within columns of the large  
204 and intermediate dataset nucleotide sequence alignments that corresponded with  
205 nucleotides identified by NASP as being paired within the HCSS. Whereas these paired  
206 nucleotides were referred to as occurring at “paired-sites”, all other nucleotides were  
207 referred to as occurring at “unpaired-sites”.

208

#### 209 **Neutrality tests for elevated negative selection at paired-sites**

210 Structural elements that increase the fitness of virus genomes are expected to be  
211 selectively preserved such that selection disfavours nucleotide substitutions should be  
212 stronger at paired-sites than at unpaired-sites. Specifically, paired-sites might display  
213 stronger evidence of purifying selection than unpaired-sites in neutrality tests such as  
214 those proposed by Tajima (70) and Fu and Li (71). We calculated Tajima’s D and Fu and  
215 Li’s F statistics for the paired- and unpaired-sites in each of the 23 large datasets.

216

217 Since in all 23 of the analysed large datasets there were invariably fewer paired-sites (i.e.  
218 those paired within the HCSS) than there were unpaired-sites (i.e. the remainder of sites in

219 the various datasets) we devised a permutation test involving the random selection of  
220 identical numbers of paired and unpaired-sites and the comparison of summary selection  
221 statistics between these paired- and unpaired-site datasets.

222

223 From each large dataset we produced 100 datasets each consisting of sites (i.e. entire  
224 large dataset alignment columns) randomly sampled with replacement from the pool of  
225 unpaired-sites. These permutation datasets contained the same numbers of sites as their  
226 corresponding paired-site datasets. Tajima's D and Fu and Li's F statistics were then  
227 calculated for all of the paired-site and permutation datasets. For each of the 23 datasets,  
228 the probability that purifying selection was operating more strongly on paired-sites than on  
229 unpaired-sites was calculated as being approximately equivalent to the proportion of times  
230 the D and F statistics calculated for the paired-site dataset were lower than those  
231 calculated for the 100 permuted datasets.

232

### 233 **Codon-based tests of synonymous substitution rates at paired versus unpaired** 234 **genomic sites.**

235 Biologically important structural elements that occur within protein-coding sequences are  
236 expected to display both selection at the codon level which favours the preservation of  
237 amino acid sequences (i.e. selection disfavours non-synonymous substitutions), and  
238 selection at the nucleotide level favouring the maintenance of base-pairing interactions  
239 within the structural elements. This "double selection" at codons that contain constituent  
240 nucleotides which form base-pairs within biologically important secondary-structures

241 should result in such codons displaying synonymous substitution rates that are lower than  
242 those occurring in codons consisting of unpaired nucleotides.

243

244 To determine whether codons corresponding to paired genomic sites displayed  
245 significantly lower synonymous substitution rates than those occurring at unpaired  
246 genomic sites, nucleotide sequences corresponding to known/suspected genes were  
247 extracted from each of the 23 intermediate dataset alignments. Within each of the resulting  
248 “gene datasets” all sites encoding amino acids in two or more different frames were  
249 removed. Following this, 43 gene datasets containing 200 or more nucleotide sites were  
250 retained for further analysis (see Table S1 for details of these datasets). Gene datasets  
251 excluded from this set because they retained too few sites included the *ren*, *mp* and *trap*  
252 genes of the GeminiTYLCV, GeminiEACMV and GeminiMYVYV datasets; ORF2 and  
253 ORF3 of the AnelloTTSuV1, AnelloTTSuV2 and AnelloTTV datasets; the *vp1* and *vp2* of  
254 the ParvoHBoV dataset; and the *vp1* of the ParvoMPV.

255

256 Two methods were used to estimate synonymous substitution rates at individual codon  
257 sites within the 43 gene datasets: Partitioning Approach for Inference of Selection  
258 (PARRIS) (72) and Fast Unconstrained Bayesian Approximation (FUBAR) (73). Both of  
259 these methods apply the time-reversible MG94 codon substitution model which utilises a  
260 61 X 61 codon substitution matrix (74) and both allow independent distributions for  
261 synonymous and non-synonymous rates. PARRIS is a random effects likelihood (REL)  
262 method permitting the use of only three discrete categories for each rate. FUBAR on the  
263 other hand is an approximate Bayesian method which permits the use of many more rate

264 classes (20 in our case) so as to increase the resolution with which, for example, subtle  
265 differences in selection pressures operating on individual codons can be distinguished.

266

267 Both FUBAR and PARRIS rely on the use of phylogenetic trees to describe the  
268 evolutionary relationships of the sequences being analysed. While it is well established  
269 that genetic recombination undermines the accuracy of phylogenetic inference (and by  
270 extension many phylogenetics-oriented codon-based selection analysis methods) (75), it  
271 was likely that many of the sequences being analysed here were recombinant (76–80). It  
272 was therefore necessary to take steps to explicitly account for recombination within these  
273 analyses. Accordingly, prior to selection analyses the Genetic Algorithm for Recombination  
274 Detection (GARD (81)) method was used to detect recombination breakpoint sites. These  
275 sites were then used to partition the alignment into “mostly recombination free” (it is  
276 unlikely that every recombination event was detected and accounted for) sub-alignments.  
277 For each of these sub-alignments a phylogenetic tree was inferred and these trees were  
278 collectively used as inputs for the PARRIS and FUBAR analyses – both of which allow  
279 phylogenetic tree topologies and branch lengths to vary across different alignment  
280 partitions so as to facilitate accurate inference of natural selection in the presence of  
281 recombination (72, 73).

282

283 Within each gene dataset, each codon was categorised as being a “paired-codon” if its  
284 third position nucleotide was paired within the relevant HCSS, and an “unpaired-codon” if  
285 its third position nucleotide was not a paired nucleotide within the relevant HCSS. Using a  
286 Mann-Whitney U test, we determined whether in each of the 43 gene datasets paired-

287 codons had significantly lower synonymous substitution rates than unpaired-codons. All  $p$ -  
288 values thus calculated were step-down corrected to account for multiple testing.

289

# 290 **Testing whether paired-sites complementarily coevolve**

291 Mutations at paired-sites may be tolerable within biologically important structural elements  
292 if they are followed by compensatory mutations that restore base-pairing (51). We  
293 detected evidence of complementary coevolution between pairs of sites within the large  
294 datasets using a customised version of the SPIDERMONKEY coevolution script written in  
295 HYPHY (82). For any chosen pair of sites within a large dataset, the script compares the  
296 standard independent sites 4 X 4 HKY85 nucleotide substitution model (83) to a 16 X 16  
297 Muse-Modified HKY85 coevolution model (called M95; (84)) to determine which of these  
298 best describes the evolution of individual site pairs. In our case entries in the M95 16 X 16  
299 substitution matrix representing changes that potentially maintain base-pairing (including  
300 both Watson-Crick pairings such as A-T and G-C and the wobble pair T-G) are multiplied  
301 by a pairing factor,  $\lambda$ , and those involving changes between paired- and unpaired-states  
302 are multiplied by  $1/\lambda$ . A maximum likelihood ratio test enabled us to determine whether  
303 nucleotide pairs were coevolving. Whereas  $\lambda > 1$  for particular coevolving site pairs  
304 indicated that they displayed a tendency towards complementary coevolution,  $\lambda = 1$   
305 indicated a tendency towards site pairs evolving independently and  $\lambda < 1$  indicated a  
306 tendency towards them coevolving non-complementarily. We identified site pairs  
307 displaying strong evidence of complementary coevolution as those with both associated  
308 maximum likelihood estimates of  $\lambda > 1$ , and Muse 95 versus HKY85 likelihood ratio test  $p$ -  
309 values  $< 0.05$ .

310

311 Importantly, due to computational intensity considerations, we restricted our analyses to  
312 testing for coevolution only between (i) pairs of sites that were within 100 nucleotides of  
313 one another and (ii) pairs of nucleotides that were polymorphic in the input dataset. Also,  
314 since recombination can undermine the accuracy with which the phylogenetic trees used  
315 to detect coevolution reflect the actual evolutionary relationships of the analysed  
316 sequences, we took steps to account for recombination in our analyses. Each large  
317 dataset was analysed for recombination using Recombination Detection Program (RDP)  
318 version 4.13 (85) which produced a “distributed alignment” in which fragments of  
319 recombinant sequences derived from different parental viruses were split up into different  
320 sequences. For each of the 23 distributed alignments obtained, a 100 nucleotide sliding  
321 window was moved 1 nucleotide step at time across the alignment. At every step the N  
322 longest nucleotide sequences were selected (where N is the number of sequences in the  
323 original alignment), and saved to an alignment file. Sequences in consecutive windows  
324 containing exactly the same N longest sequences were merged into one file (ensuring that  
325 no sites from the original alignment were duplicated in the merged alignment). Maximum-  
326 likelihood (ML) phylogenetic trees were inferred for each of the resulting alignments under  
327 the HKY85 nucleotide substitution model using PhyML3.0 (86). Each of these alignments  
328 and their corresponding phylogenetic tree was used as inputs for our complementary  
329 coevolution analysis.

330

### 331 **Customised computational tools**

332 All computer scripts used in all the analyses conducted can be downloaded from  
 333 <http://web.cbio.uct.ac.za/~brejnev/downloads/Scripts.zip> and a customised computer  
 334 program and datasets for visualisation of predicted structural elements can be downloaded  
 335 from <http://web.cbio.uct.ac.za/~brejnev/downloads/DOOSS.zip> (unzip these files and  
 336 please see the README file for instructions).

337

## 338 **Results and Discussion**

### 339 **Numerous evolutionarily conserved secondary-structures are evident within** 340 **eukaryotic ssDNA virus genomes**

341 We assembled 23 full-genome sequence datasets representing the families *Circoviridae*,  
 342 *Anelloviridae*, *Parvoviridae*, *Nanoviridae* and *Geminiviridae* (Table 1). In each of these  
 343 between 69 and 316 conserved secondary-structural elements were identified using the  
 344 minimum free energy (MFE) approach implemented in the computer program NASP (66).  
 345 From these lists of conserved structural elements, NASP identified subsets of between five  
 346 and 132 “high confidence” structural elements. These lists, hereafter referred to as “high  
 347 confidence structure set” (HCSS) lists, contained those structures primarily responsible for  
 348 the analysed genomes having greater degrees of predicted structural stability than those  
 349 of randomised sequences with identical nucleotide compositions. Notably, most of the  
 350 previously described biologically important structures in these viral genomes were present  
 351 within the top 30% of structures in the HCSS lists of their respective datasets. These  
 352 included hairpin structures at the virion strand origins of replication in circoviruses,  
 353 nanoviruses, and geminiviruses and T-shaped structures required for replication in



354 parvoviruses. The genomic coordinates of structures within all 23 of the HCSSs were  
355 mapped onto their respective genomes (Fig. 1; Fig. 2; Table S2).

356

357 Clearly our computational approach for predicting secondary-structure suggests that there  
358 exist more conserved genomic secondary-structures within many of these ssDNA virus  
359 genomes than is currently appreciated. It is plausible that, as is the case with currently  
360 known secondary-structures within these genomes, many of the uncharacterised  
361 conserved structures may have been preserved during evolution due to their biological  
362 importance.

363

364 Although directly testing the biological relevance of any one of the identified potential  
365 structures would require detailed mutational analyses of their constituent nucleotides  
366 within the context of infectious cloned genomes or analysis of recombinant viral  
367 constructs, followed by extensive quantitative fitness assays (87, 88), there are less  
368 cumbersome computational approaches for testing whether the identified structures  
369 collectively (as opposed to individually) are likely to have any biological relevance. In this  
370 regard the biological relevance of the structures in our HCSSs could be tested by  
371 comparing how their constituent nucleotides evolve relative to those at positions located  
372 outside of the HCSSs.

373

374 Therefore, in our subsequent analyses we focused on testing whether, relative to the  
375 remainder of nucleotides in the genomes (hereafter referred to as unpaired nucleotides),  
376 nucleotides predicted to be base-paired within the HCSS (hereafter referred to as paired

377 nucleotides), are evolving in ways suggestive of their parental structural elements  
378 possessing some biological function.

379

# 380 **Purifying selection is apparently strongest at paired nucleotide sites**

381 Nucleotide sites involved in biologically important base-pairing interactions might be  
382 expected to evolve under a greater degree of purifying selection (selection against  
383 change) than sites that are not base-paired. Also, sequences evolving under purifying  
384 selection are expected to have lower frequencies of minor allele polymorphisms than  
385 those evolving under neutral selection and are, therefore, expected to yield negative  
386 values for Tajima's D and Fu and Li's F statistics (70, 71). If purifying selection was  
387 stronger at base-paired-sites than at unpaired-sites we would expect to see lower values  
388 of the D and F statistics for datasets containing only base-paired-sites (constructed from  
389 the large dataset alignments by removing all unpaired nucleotide sites) than for datasets  
390 containing only unpaired-sites (constructed from the large dataset alignments by removing  
391 all paired nucleotide sites).

392

393 In all but one of the 23 large datasets (the exception being the circovirus dataset,  
394 CircoDGCV; Table 1), both the paired and unpaired-site alignments consistently yielded  
395 negative D and F test static values (see Table 2). In 16/23 of the datasets both the D and  
396 F statistics were lower for the paired-site than for the unpaired-site datasets. In 5/23  
397 datasets (the anellovirus datasets AnelloTTSuV2 and AnelloTTV, the parvovirus datasets,  
398 ParvoAAV, ParvoHBoV and the geminivirus dataset GeminiTYLCV) either the D or F  
399 statistics were lower for the paired-site datasets than for unpaired-site datasets. In only

400 2/23 cases (CircoDGCV and AnelloTTSuV1) did the unpaired-site dataset yield both  
401 values of D and F statistics lower than those yielded by the paired-site datasets.

402

403 This observation is consistent with our hypothesis that, if paired-sites within the 23 HCSS  
404 lists really do reside within biologically important secondary-structures, they should display  
405 higher degrees of purifying selection than other sites within the analysed genomes.

406

407 However, to test whether values of these statistics were significantly lower at paired-sites  
408 than at unpaired-sites in the 21/23 large datasets displaying the expected trend, for each  
409 dataset we applied a permutation test involving resampling of identical numbers of sites  
410 from the unpaired-site dataset as were present within the paired-site dataset (in each  
411 dataset unpaired-sites were invariably more numerous than the paired-sites). In each case a  $p$ -  
412 value was computed as the proportion of the 100 permuted unpaired-site datasets that  
413 yielded lower D or F values than the corresponding paired-site dataset. In this test a  $p$ -  
414 value  $< 0.05$  indicates that you would expect to see a D or F value for an unpaired-site  
415 dataset that was lower than that of its corresponding paired-site dataset less than 5% of  
416 the time if the null model of neutral evolution was true.

417

418 In 11/23 datasets both the D and F statistic permutation tests yielded evidence that paired-  
419 sites within the HCSS lists experience significantly stronger ( $p$ -values  $< 0.05$ ) purifying  
420 selection than the remainder of genomic sites. In a further 6/23 cases, either the D, or F  
421 statistic test yielded at least marginal evidence ( $p$ -values  $< 0.08$ ) of paired-sites  
422 experiencing stronger purifying selection than unpaired-sites. Therefore, in only 6/23

423 cases, was there absolutely no evidence of paired-sites experiencing significantly stronger  
424 purifying selection than unpaired-sites.

425

426 Interestingly, all three of the analysed anellovirus datasets were among the six datasets  
427 with no evidence of purifying selection acting on paired-sites. It is perhaps also noteworthy  
428 that of the eleven datasets displaying strong evidence of base-pairing associated with  
429 purifying selection, only two (both of them circoviruses, CircoCoCV and CircoBFDV ) were  
430 from the ten mammal- and bird-infecting virus datasets. While it is not possible to directly  
431 compare the plant- and animal-infecting virus datasets to one another, it is plausible that  
432 increased structural stability afforded by the lower physiological temperatures of plants  
433 relative to animals might contribute to the genomic structures of the plant viruses being  
434 more evolutionarily stable than those of their warm-blooded animal counterparts. A more  
435 mundane explanation, however, could simply be that our animal virus datasets were, in  
436 general, substantially smaller than our plant virus datasets and that our analysis therefore,  
437 simply lacked sufficient power to differentiate between the numbers of low frequency  
438 polymorphisms within the paired and unpaired dataset fractions.

439

440 Regardless of possible differences between animal and plant viruses, collectively these  
441 results indicate that a substantial proportion of paired-sites within at least 17/23 of the  
442 HCSSs are evolving in a manner that is consistent with many of these structures being  
443 evolutionarily preserved.

444

445 **Synonymous substitution rates are unusually low at paired genomic sites**

446 We hypothesised that selection favouring the maintenance of base-pairing within  
 447 secondary-structures might be particularly evident when these structures occurred within  
 448 protein-coding regions of the genome. Essentially, we investigated whether codons in  
 449 which third codon position nucleotides were predicted to be base-paired within the HCSSs  
 450 had significantly lower synonymous substitution rates than those with unpaired nucleotides  
 451 in the third codon position.

452 Synonymous substitution rates at individual codon sites within 43 gene datasets (Table  
 453 S1) were inferred using the random effects likelihood selection analysis methods, PARRIS  
 454 (75) and FUBAR (73). These methods indicated that in 27/43 of these datasets, the  
 455 median substitution rates of codons with paired third position nucleotides were significantly  
 456 lower than those of codons with unpaired third position nucleotides (multiple comparison  
 457 corrected Mann-Witney U-test  $p$ -value  $< 0.05$ ). An additional five datasets yielded similar  
 458 evidence but only with one of the two selection analysis methods (Table 3).

459 The results of these analyses therefore strongly support our hypothesis that two layers of  
 460 selection – one operating at the amino acid sequence level and the other at the nucleotide  
 461 sequence level – are likely acting on nucleotide sites within the HCSSs that fall within  
 462 coding regions. This suggests not only that many of the predicted secondary-structures  
 463 represented within the HCSSs really do exist (either within single-stranded genomic DNAs  
 464 themselves, or within the RNA transcripts that are produced from them), but that these  
 465 structures likely also make a substantial contribution to the fitness of the genomes within  
 466 which they reside.

467

468 While evidence of lower degrees of nucleotide polymorphism and decreased synonymous  
 469 substitution rates at paired-sites than at unpaired-sites provides strong support for the  
 470 existence of many of the predicted secondary-structural elements within the HCSSs, it  
 471 must be stressed that this result does not necessarily imply that these elements are  
 472 biologically functional. The reason for this is that besides influencing which arising  
 473 mutations are deleterious and which are neutral (and, therefore, which mutations are likely  
 474 to be purged from populations by natural selection), the presence of secondary-structures  
 475 within ssDNA genomes could potentially also influence the basal rates at which sites  
 476 within these genomes become mutated (89), simply because base-paired nucleotides  
 477 might be predisposed to lower mutation rates than their unpaired counterparts (90, 91).

478

479 **In short term evolution experiments mutations tend to preferentially accumulate at**  
 480 **unpaired-sites**

481 If paired-sites within the HCSSs really do form base-pairs within genomic secondary-  
 482 structures, we hypothesized that these sites might accumulate fewer mutations than  
 483 unpaired-sites. We tested this hypothesis using mutation data from a series of previously  
 484 published short-term evolution experiments. In one experiment infectious cloned genomes  
 485 of two *Maize streak virus* isolates (called MSV-MatA and MSV-VW) closely related to  
 486 those in the GeminiMSV dataset, were used to infect maize plants (92). In another  
 487 experiment infectious cloned genomes of a *Tomato yellow leaf curl virus* isolate (called  
 488 TYX) and a *Tomato leaf curl Comoros virus* isolate (called TOX; both closely related to  
 489 sequences included in the GeminiTYLCV dataset) were used to infect tomato plants (52).

490

491 While over 101 days post-infection the MSV-MatA and MSV-VW genomes were noted to  
492 have accumulated 41 and 33 mutations, respectively, at 52 distinct nucleotide sites, over  
493 120 days the TYX and TOX genomes had respectively accumulated 31 and 105 mutations  
494 at 135 distinct nucleotide sites. As described previously for our small datasets, we  
495 predicted the secondary-structures of each genome pair using NASP in order to obtain, for  
496 each pair, its own specific HCSS. We used these HCSSs to construct two-by-two  
497 contingency tables for paired-sites (sites predicted to be paired within the HCSS) and  
498 unpaired-sites (all sites in the genome other than the HCSS paired-sites) versus variable  
499 sites (those where mutations occurred) and invariable sites (those where mutations did not  
500 occur) and used these in a Fisher's exact test (93), to assess whether variable sites were  
501 significantly clustered outside rather than inside paired-sites.

502

503 For MSV-MatA and MSV-VW 11/52 variable sites (~21%) were located at paired  
504 nucleotide sites (939/2641 or ~36% of considered sites) within the HCSS, yielding  
505 significant evidence ( $p$ -value = 0.019) that mutations tended to occur more frequently at  
506 unpaired nucleotides. Similarly, for TYX and TOX only 5/135 variable sites (~4%) were  
507 located at paired nucleotide sites (237/2724 or ~9% of considered sites) within the HCSS  
508 regions, indicating a significant tendency ( $p$ -value = 0.021) for mutations to accumulate  
509 more frequently at unpaired nucleotide sites.

510

511 Although no analogous experimental data is currently available for any of the other plant-  
512 and animal-infecting ssDNA viruses investigated here, it is nevertheless important that  
513 even in short term geminivirus evolution experiments such as these, where selection has

514 not had prolonged periods to purge slightly deleterious mutations, there remains such an  
515 obvious trend for mutations to preferentially occur at unpaired-sites.

516

517 Unfortunately, even though these experiments were short-term (lasting between 101 and  
518 120 days), it remains possible that selection, in addition to a decreased biochemical  
519 predisposition to mutation, was responsible for the relatively lower mutation frequencies at  
520 paired-sites within these genomes. While still consistent with our hypothesis that selection  
521 is acting on secondary-structures to maintain their biological functionality, these results  
522 suggest that the alternative hypothesis – that base-paired-sites within secondary-  
523 structures are simply biochemically predisposed to mutate more slowly than unpaired-sites  
524 – is also entirely plausible.

525

526 Therefore, although we had established up to this point that secondary-structures are  
527 likely quite pervasive within ssDNA virus genomes, we were unable to definitively attribute  
528 the apparent evolutionary conservation of these structures to natural selection favouring  
529 the maintenance of their biological functionality.

530

#### 531 **Base-paired-sites tend to complementarily coevolve**

532 It is expected that, independent of different basal mutation rates at paired- and unpaired-  
533 sites, nucleotide substitutions that occur at paired-sites within biologically functional  
534 secondary-structures might only be tolerable if coupled with complementary substitutions  
535 that reconstitute base-pairing. Therefore, in order to test for natural selection acting to  
536 maintain secondary-structures without the confounding effects of base-pairing dependent



537 basal mutation rate variation, we directly tested for evidence of paired-sites within the  
 538 HCSSs coevolving with one another in a manner consistent with the maintenance of their  
 539 base-pairing. Specifically, we tested for associations between sites predicted to be base-  
 540 paired within the HCSSs and sites detectably coevolving in a complementary fashion  
 541 within the 23 large datasets. For each large dataset we performed a two-by-two  
 542 contingency test of site pairs predicted to be paired versus unpaired on the one hand, and  
 543 sites predicted to be coevolving versus not coevolving on the other.

544

545 In all but one circovirus dataset, CircoCoCV, we found strong significant associations  
 546 (multiple testing corrected  $p$ -values  $<0.0001$ ) between paired-sites within the HCSSs and  
 547 sites for which complementary coevolution was detected (Table 4). It is noteworthy that  
 548 the CircoCoCV was one of the two animal-infecting virus datasets displaying both strong  
 549 evidence of base-pairing associated negative selection, and evidence of strong selection  
 550 disfavouring synonymous substitutions at paired codon sites within coding regions.  
 551 Therefore, the lack of significant evidence of coevolution between nucleotides predicted to  
 552 be paired within the CircoCoCV HCSS may simply be due to strong selection disfavouring  
 553 any substitutions at these sites.

554

555 Besides providing additional evidence that many of the structures represented within the  
 556 HCSSs really do form either within the genomes of these ssDNA viruses, or within their  
 557 RNA transcripts, this result provides the most compelling evidence yet that natural  
 558 selection is favouring the maintenance of a substantial proportion of these structures. The  
 559 simple fact that many of the structures represented within the HCSSs likely provide

560 significant fitness advantages to the genomes in which they occur, in turn, suggests that  
561 many of these structures have as yet, undetermined biological functions.

562

### 563 **Potentially important structural elements within eukaryotic ssDNA virus genomes**

564 Whereas we provided evidence of pervasive evolutionarily conserved (and therefore, likely  
565 biologically functional) secondary-structures within the various ssDNA virus genomes that  
566 we have analysed, we have not up to this point examined any of the individual  
567 computationally inferred structural elements in any significant detail. Fortunately, some of  
568 the analyses that we performed provide a straightforward means of ranking the identified  
569 structures within the HCSSs in order of their likely biological functionality (94). Specifically,  
570 these rankings were based on: (1) the degree to which structural elements were  
571 conserved across the analysed genomes; (2) the degree to which synonymous  
572 substitution rates were constrained at codon sites containing nucleotides that are  
573 predicted to be base-paired and; (3) the degree to which nucleotides predicted to be base-  
574 paired coevolve with one another. Rankings based collectively on these three criteria are  
575 hereafter, referred to as “consensus rankings” (Table S2).

576

577 The ten highest ranked structural elements based on these criteria within each of the 23  
578 analysed HCSSs are plotted in magenta and cyan in Fig. 1 and Fig. 2, and are listed in  
579 Table S2. It is important to point out that although these top ranked structures contributed  
580 most to the signals detected in our earlier association tests, it is possible that many of  
581 them do not actually exist in the exact form that we have inferred either in the ssDNA  
582 genomes themselves, or in the RNA molecules transcribed from these genomes. Besides

583 expected inaccuracies in the computational inference of DNA and RNA secondary-  
 584 structures (95), it is likely that even if these structural elements have been accurately  
 585 inferred, the exact base-pairing configurations within the presented consensus structures  
 586 will likely vary between the different genomes within each of the analysed datasets. Also, it  
 587 is very likely that, even within an individual genome, many of these structures will not be  
 588 static but will instead represent a single reasonably stable base-pairing configuration  
 589 amongst a (potentially very large) ensemble of similarly stable alternative configurations. It  
 590 should therefore, be borne in mind that the actual base-pairing interactions within the  
 591 tertiary structures represented by many of these structural elements, might vary as the  
 592 structural elements continually transition between their alternative forms.

593

594 Among the individual structural elements that achieved the highest consensus rankings  
 595 were all of the well-characterised secondary-structures found at the origins of replication of  
 596 circoviruses (ranks 1 to 6), nanoviruses (ranks 8 to 28), geminiviruses (ranks 1 to 12) and  
 597 parvoviruses (ranks 1 to 35) (Table S2).

598

599 Additional well-characterised structures detected include the replication associated protein  
 600 gene (*rep*) intron associated structure (GeminiMSV; rank 16) (51), the parvovirus  
 601 transcription attenuation stem-loop structures (ParvoMPV; ranks 17 and 34; Table S2)  
 602 (42), the 3' complementary strand T-shaped structure that binds to the viral capsid in some  
 603 parvoviruses (ParvoMPV; rank 3; Table S2) (96).

604

Besides these well-known structures, we sought to identify other uncharacterised, but likely biologically functional, structural elements within some of these genomes. Rather than exhaustively enumerating every predicted secondary-structural element that might have some biological relevance, we instead focus here on a few examples of the elements that have apparently been conserved across multiple, highly divergent viral lineages in the various viral families that we analysed.

## Geminivirus

We identified a particularly conserved 126 to 157 nt secondary-structure within the movement protein (*mp*) gene of all five analysed mastrevirus datasets (GeminiMSV, GeminiPanSV, GeminiWDV, GeminiTYDV-CpCV and GeminiCpCDV; structure G1 in Fig. 1 and Fig. 3). In all of these datasets other than GeminiMSV, the entire structure was within the HCSS (7<sup>th</sup> out of 25 in GeminiCpCDV, 10<sup>th</sup> out of 47 in GeminiTYDV-CpCV, 21<sup>st</sup> out of 31 in GeminiPanSV, 16<sup>th</sup> out of 30 in GeminiWDV, and 51<sup>st</sup> in GeminiMSV). The structure in the GeminiMSV dataset displayed a particularly high degree of conformational similarity with that in the GeminiPanSV dataset with the two structures sharing a nearly identical 21-nucleotide long stem sequence (Fig. 3) indicating that they are almost certainly homologous. Although the sequences within this structure differ substantially between the other mastrevirus datasets, they all contain the splice donor, acceptor and branch sites previously identified (or predicted) in mastrevirus *mp* introns (97) (Fig. 3), suggesting that the structure is possibly functional within the *mp* mRNA transcript where it might facilitate *mp* intron splicing. Also, likely acceptor and donor sites identified within these various sequences tend to occur at junctions between paired and unpaired

628 nucleotides – a factor which might enhance the accessibility of these sites during splicing  
629 (18, 98, 99).

630

631 Another highly conserved secondary-structure that is most likely functional within  
632 geminivirus genomes was identified near the 3' end of the coat protein (*cp*) genes of  
633 begomoviruses in the GeminiTYLCV, GeminiEACMV and GeminiMYVYV datasets  
634 (structure G2 in Fig. 1 and Fig. 4). This structure contains a conserved stem-loop  
635 sequence immediately 3' of the *cp* stop codon that contains the likely polyadenylation  
636 signals of both virion and complementary strand RNA transcripts (Fig. 4). It is likely  
637 therefore, that this structure may be functional either within ssDNA as a transcriptional  
638 terminator, or within transcribed mRNA during polyadenylation.

639

640 Parvovirus

641 We identified a variety of uncharacterised parvovirus genomic and/or mRNA structural  
642 elements with potential functionality at the start of the large non-structural (*ns1*) gene  
643 (Structure P1 in Fig. 2 and Fig. 5; 20<sup>th</sup> out of 70 HCSS structures in ParvoAAV, 30<sup>th</sup> out of  
644 105 HCSS structures in ParvoHBoV and 34<sup>th</sup> out of 132 HCSS structures in ParvoMPV),  
645 the start of the major virion/viral protein (*vp1*) gene (Structure P2 in Fig. 2 and Fig. 5; 16<sup>th</sup>  
646 out of 70 HCSS structures in ParvoAAV and 9<sup>th</sup> out of 132 HCSS structures in ParvoMPV),  
647 the start of the small non-structural (*np1*) gene (structure P3 in Fig. 2 and Fig. 5, 59<sup>th</sup> out of  
648 105 HCSS structure in ParvoHBoV) and the start of the minor virion protein (*vp2*) gene  
649 (structure P4 in Fig. 2 and Fig. 5; 56<sup>th</sup> out of 132 HCSS structures in ParvoMPV). Although  
650 there were no sequence similarities shared between positionally analogous structures in

651 the different parvovirus datasets, this was not unexpected given that these datasets  
 652 represent species within different genera (with sequences in different datasets sharing on  
 653 average only 57.8% sequence identity). The ParvoMPV IR-*ns1* structure contains a stem-  
 654 loop identified to play role in transcription attenuation of *Parvovirus minute virus of mice*  
 655 (42) (structure P1; Fig. 5). In this regard, it is noteworthy that start codons within the  
 656 structures that we have identified are consistently located either within, or immediately  
 657 adjacent to unpaired loop or bulge regions (Fig. 5). This tendency was also noted in other  
 658 datasets analysed, and it is plausible that structures spanning the start codons of genes in  
 659 these different families are functional within either partially single stranded DNA during the  
 660 initiation of transcription, or in transcribed mRNA during the initiation of translation.

661

#### 662 Circovirus

663 While we were unable to identify any secondary-structures that were clearly conserved  
 664 across all five circovirus datasets analysed, within the moderately divergent CircoCoCV  
 665 and CircoBFDV datasets (these two datasets share on average 64% pairwise sequence  
 666 identity), we identified an intergenic region (IR) stem-loop structure (structure C1 in Fig. 2  
 667 and Fig. 6), which is highly conserved in each of the respective datasets (ranked 5<sup>th</sup> out of  
 668 35 HCSS structures in CircoCoCV and 1<sup>st</sup> out 41 HCSS structures in CircoBFDV). Despite  
 669 the sequences of this structure sharing no obvious similarity between the two datasets, in  
 670 both datasets the stem is GC rich (and therefore, predicted to be very stable) with a loop  
 671 sequence containing a conserved pentanucleotide (CGAAG). This structural element  
 672 could potentially contain the complementary strand replication origin or it might be

673 functional either during the termination of transcription, or in the post-transcriptional  
674 processing of mRNA transcripts.

675

676 Anellovirus

677 A conserved T-shaped structure was identified in the IRs of the two anellovirus datasets,  
678 AnelloTTSuV1 and AnelloTTuSV2 (structure A1 in Fig. 2 and Fig. 7; 7<sup>th</sup> out of 20  
679 structures in the AnelloTTSuV1 HCSS and 1<sup>st</sup> out of 27 structures in the AnelloTTSuV2  
680 HCSS). Even though these two datasets are moderately divergent (sequences within them  
681 share on average 60.7% pairwise identity), the structure is strikingly conserved between  
682 the two datasets. In both datasets it has a nearly identical predicted T-shaped  
683 conformation with a highly conserved 17-nucleotide long sequence at the top of the “T”  
684 (highlighted in sky blue in Fig. 7).

685

686 Given the high degree to which this structure has been conserved between these two  
687 moderately divergent anellovirus datasets, we attempted to identify a homologous  
688 structure within our third highly divergent anellovirus dataset (AnelloTTV). The most likely  
689 homologue of this structure also resides within the IR and is ranked 7<sup>th</sup> out of 78 structures  
690 in the AnelloTTV HCSS (structure A1 in Fig. 2 and Fig. 7). However, the AnelloTTV  
691 structure has a stem-loop rather than a T-shaped configuration and lacks the 17-  
692 nucleotide sequence that is conserved in the AnelloTTSuV1 and AnelloTTSuV2 structures.  
693 All three Anellovirus structures nevertheless, contain two similar sequences (five- and six  
694 nucleotide long) at similar positions within their stems (outlined in blue and yellow in Fig.  
695 7), which strongly suggests that these structures are indeed homologous.

696

697 Unlike with many other circular ssDNA viruses that replicate by rolling circle replication, it  
698 is presently unknown where the Anellovirus virion and complementary strand origins of  
699 replication (*oris*) reside. Given that the virion strand *oris* of other ssDNA viruses generally  
700 occur within IRs and have a characteristic stem-loop structure with an A-T rich loop  
701 sequence, it is plausible that this highly conserved Anellovirus structural element might  
702 contain the Anellovirus virion strand *ori*. However, characterisation of replication  
703 competent sub-full-length *Torque teno virus* (TTV) genomes (which are closely related to  
704 those represented in our AnelloTTV dataset) has suggested that the TTV virion strand *ori*  
705 is approximately 470 nucleotides 3' of the highly conserved TTV stem-loop structure that  
706 we have identified here (in the region of a small stem-loop structure ranked 83<sup>rd</sup>, below  
707 the HCSS in our AnelloTTV dataset; data not shown) (100). Importantly, the structure we  
708 have identified falls outside the genomic region that is conserved within these defective  
709 genomes and, in the TTV genome at least, is therefore, unlikely to be the virion strand *ori*.  
710 Apart from possibly containing the virion strand *ori*, this highly conserved structural  
711 element could alternatively be involved in either complementary strand replication or  
712 transcriptional regulation, both of which are also carried out by IR sequences in all other  
713 known ssDNA viruses.

714

## 715 **Conclusion**

716 Using computational methods we have identified numerous secondary-structures that  
717 probably form at least transiently, within eukaryotic ssDNA virus genomes, and shown that  
718 a significant proportion of these predicted structures are likely biologically functional (Table



719 5). We have further provided a few examples of currently uncharacterised genomic  
720 secondary-structures which, due to high degrees of evolutionary conservation across  
721 multiple highly divergent viral lineages, likely play a central role in the biology of the  
722 various ssDNA viruses examined here.

723

724 Although we found evidence consistent with natural selection strongly disfavoured the  
725 accumulation of substitutions at paired-sites, we also found that paired-sites tended to  
726 display lower nucleotide variability than unpaired-sites. Using data from published  
727 evolution experiments, we showed that, in at least one of the analysed virus families (the  
728 geminiviruses), it is possible that this discrepancy may simply be due to mutation  
729 frequencies at paired-sites being lower than those at unpaired-sites (possibly due to base-  
730 paired nucleotides being less mutable than unpaired nucleotides). We were nevertheless  
731 able to clearly demonstrate the action of selection by showing that those base-paired-sites  
732 which do accumulate mutations display a significant tendency towards complementary  
733 coevolution with their predicted pairing partners – presumably to maintain the biological  
734 function of their parent structures.

735

736 Despite providing compelling evidence of pervasive biologically functional secondary-  
737 structures within eukaryote-infecting ssDNA viruses, it is important to reiterate that our  
738 study has certain limitations. It is very likely that the complex genomic structures of these  
739 viruses are not entirely static. The secondary and tertiary structures of these entire  
740 genomes are, in fact, very likely to shift continually between large numbers of different  
741 thermodynamically stable states. We cannot therefore, be absolutely certain if the

742 computationally predicted structures identified here are a good reflection of those which  
743 form most commonly within these ssDNA virus genomes. Also, although examples of  
744 individual genomic structural elements that are highly conserved across divergent virus  
745 lineages are likely to have some biological functionality, we cannot know without further  
746 laboratory experimentation either what the precise functions of these structures might be,  
747 or whether they function within the context of ssDNA or transcribed RNA.

748

749 Regardless of whether specific individual structures form, or are functional within ssDNA  
750 or transcribed RNA molecules, it is absolutely clear from our study that, at the whole-  
751 genome scale, selection favouring the overall maintenance of pervasive biologically  
752 functional nucleic acid secondary-structures has likely been a major theme in the  
753 evolutionary history of eukaryotic ssDNA viruses.

754

755 **Acknowledgments**

756 BMM is funded by the University of Cape Town, South Africa. PL and JML are funded by  
757 the Conseil Régional de La Réunion, European Union (FEDER) and Centre de  
758 Coopération Internationale en Recherche Agronomique pour le Développement (CIRAD).  
759 ALM is funded by the National Research Foundation (South Africa) and the Carnegie  
760 Corporation of New York. BM is funded by the CFAR Translational Virology Core: P30  
761 AI036214, Molecular Epidemiology (Avant Garde grant): DP1 DA034978. AFYP is  
762 supported by a New Investigator Award from the Canadian Institutes of Health Research  
763 (Canadian HIV Vaccine Initiative) and by a Scholar Award from the Michael Smith  
764 Foundation for Health Research / St. Paul's Hospital Foundation - Providence Health Care  
765 Research Institute. DNS is funded by the Pannar (Pty) Ltd. DPM and GWH are funded by  
766 the South African National Research Foundation. The authors would like to thank the  
767 Centre for High Performance Computing in Cape Town and the Information  
768 Communication Technology Services Department at the University of Cape Town for use  
769 of their high-performance computing clusters.

770

## 771 References

- 772 1. **Yuen L, Moss B.** 1987. Oligonucleotide sequence signaling transcriptional  
773 termination of vaccinia virus early genes. *Proc. Natl. Acad. Sci. U. S. A.* **84**:6417–21.
- 774 2. **Hefferon KL, Moon Y-S, Fan Y.** 2006. Multi-tasking of nonstructural gene products  
775 is required for bean yellow dwarf geminivirus transcriptional regulation. *FEBS J.*  
776 **273**:4482–94.
- 777 3. **Shen R, Miller WA.** 2004. The 3' untranslated region of tobacco necrosis virus RNA  
778 contains a barley yellow dwarf virus-like cap-independent translation element. *J.*  
779 *Viol.* **78**:4655–64.
- 780 4. **Song SI, Miller WA.** 2004. cis and trans Requirements for Rolling Circle Replication  
781 of a Satellite RNA. *J. Virol.* **78**:3072–3082.
- 782 5. **Stockley PG, Twarock R, Bakker SE, Barker AM, Borodavka A, Dykeman E,**  
783 **Ford RJ, Pearson AR, Phillips SE V, Ranson N a, Tuma R.** 2013. Packaging  
784 signals in single-stranded RNA viruses: nature's alternative to a purely electrostatic  
785 assembly mechanism. *J. Biol. Phys.* **39**:277–87.
- 786 6. **Steinfeldt T, Finsterbusch T, Mankertz A.** 2001. Rep and Rep' protein of porcine  
787 circovirus type 1 bind to the origin of replication in vitro. *Virology* **291**:152–60.
- 788 7. **Berns KI.** 1990. Parvovirus replication. *Microbiol. Rev.* **54**:316–29.
- 789 8. **Gronenborn B.** 2004. Nanoviruses: genome organisation and protein function. *Vet.*  
790 *Microbiol.* **98**:103–109.
- 791 9. **Ashktorab H, Srivastava A.** 1989. Identification of nuclear proteins that specifically  
792 interact with adeno-associated virus type 2 inverted terminal repeat hairpin DNA. *J.*  
793 *Viol.* **63**:3034–9.
- 794 10. **Faurez F, Dory D, Grasland B, Jestin A.** 2009. Replication of porcine circoviruses.  
795 *Viol. J.* **6**:60.
- 796 11. **Sun X, Simon AE.** 2006. A cis-replication element functions in both orientations to  
797 enhance replication of Turnip crinkle virus. *Virology* **352**:39–51.
- 798 12. **Mohan BR, Dinesh-Kumar SP, Miller WA.** 1995. Genes and cis-acting sequences  
799 involved in replication of barley yellow dwarf virus-PAV RNA. *Virology* **212**:186–95.
- 800 13. **Ilyinskii PO, Schmidt T, Lukashev D, Meriin AB, Thoidis G, Frishman D,**  
801 **Shneider AM.** 2009. Importance of mRNA secondary structural elements for the  
802 expression of influenza virus genes. *OMICS* **13**:421–30.

- 803 14. **Koev G, Mohan BR, Miller WA.** 1999. Primary and secondary structural elements  
804 required for synthesis of barley yellow dwarf virus subgenomic RNA1. *J. Virol.*  
805 **73**:2876–85.
- 806 15. **Guo L, Allen EM, Miller WA.** 2001. Base-pairing between untranslated regions  
807 facilitates translation of uncapped, nonpolyadenylated viral RNA. *Mol. Cell* **7**:1103–  
808 9.
- 809 16. **Zuo X, Wang J, Yu P, Eyler D, Xu H, Starich MR, Tiede DM, Simon AE.** 2009.  
810 Solution structure of the cap-independent translational enhancer and ribosome-  
811 binding element in the 3' UTR of turnip crinkle virus. *Biophys. Comput. Biol.*  
812 **107**:1385–1390.
- 813 17. **Watts JM, Dang KK, Gorelick RJ, Leonard CW, Bess JW, Swanstrom R, Burch**  
814 **CL, Weeks KM.** 2009. Architecture and secondary structure of an entire HIV-1 RNA  
815 genome. *Nature* **460**:711–6.
- 816 18. **Moss WN, Dela-Moss LI, Priore SF, Turner DH.** 2012. The influenza A segment 7  
817 mRNA 3' splice site pseudoknot/hairpin family. *RNA Biol.* **9**:1305–10.
- 818 19. **Wikström FH, Meehan BM, Berg M, Timmusk S, Elving J, Fuxler L, Magnusson**  
819 **M, Allan GM, McNeilly F, Fossum C.** 2007. Structure-dependent modulation of  
820 alpha interferon production by porcine circovirus 2 oligodeoxyribonucleotide and  
821 CpG DNAs in porcine peripheral blood mononuclear cells. *J. Virol.* **81**:4919–27.
- 822 20. **Wikström FH, Fossum C, Fuxler L, Kruse R, Lövgren T.** 2011. Cytokine induction  
823 by immunostimulatory DNA in porcine PBMC is impaired by a hairpin forming  
824 sequence motif from the genome of Porcine Circovirus type 2 (PCV2). *Vet. Immunol.*  
825 *Immunopathol.* **139**:156–66.
- 826 21. **Simmonds P, Tuplin A, Evans DJ.** 2004. Detection of genome-scale ordered RNA  
827 structure (GORS) in genomes of positive-stranded RNA viruses: Implications for  
828 virus evolution and host persistence. *RNA* **10**:1337–51.
- 829 22. **Simon AE, Gehrke L.** 2009. RNA conformational changes in the life cycles of RNA  
830 viruses, viroids, and virus-associated RNAs. *Biochim. Biophys. Acta* **1789**:571–83.
- 831 23. **Fernandes J, Jayaraman B, Frankel A.** 2012. The HIV-1 Rev response element:  
832 an RNA scaffold that directs the cooperative assembly of a homo-oligomeric  
833 ribonucleoprotein complex. *RNA Biol.* **9**:6–11.
- 834 24. **Powell DM, Amaral MC, Wu JY, Maniatis T, Greene WC.** 1997. HIV Rev-  
835 dependent binding of SF2/ASF to the Rev response element: possible role in Rev-  
836 mediated inhibition of HIV RNA splicing. *Proc. Natl. Acad. Sci. U. S. A.* **94**:973–8.

- 837 25. **Le SY, Malim MH, Cullen BR, Maizel J V.** 1990. A highly conserved RNA folding  
838 region coincident with the Rev response element of primate immunodeficiency  
839 viruses. *Nucleic Acids Res.* **18**:1613–23.
- 840 26. **You S, Stump DD, Branch AD, Rice CM.** 2004. A cis-Acting Replication Element in  
841 the Sequence Encoding the NS5B RNA-Dependent RNA Polymerase Is Required  
842 for Hepatitis C Virus RNA Replication. *J. Virol.* **78**:1352–1366.
- 843 27. **Koev G, Liu S, Beckett R, Miller WA.** 2002. The 3'-Terminal Structure Required for  
844 Replication of Barley Yellow Dwarf Virus RNA Contains an Embedded 3' End.  
845 *Virology* **292**:114–126.
- 846 28. **Raman S, Bouma P, Williams GD, Brian DA.** 2003. Stem-loop III in the 5'  
847 untranslated region is a cis-acting element in bovine coronavirus defective  
848 interfering RNA replication. *J. Virol.* **77**:6720–30.
- 849 29. **Pillai-Nair N, Kim K-H, Hemenway C.** 2003. Cis-acting Regulatory Elements in the  
850 Potato Virus X 3' Non-translated Region Differentially Affect Minus-strand and Plus-  
851 strand RNA Accumulation. *J. Mol. Biol.* **326**:701–720.
- 852 30. **Chen D, Barros M, Spencer E, Patton JT.** 2001. Features of the 3'-consensus  
853 sequence of rotavirus mRNAs critical to minus strand synthesis. *Virology* **282**:221–9.
- 854 31. **Paul A V, Rieder E, Kim DW, van Boom JH, Wimmer E.** 2000. Identification of an  
855 RNA hairpin in poliovirus RNA that serves as the primary template in the in vitro  
856 uridylation of VPg. *J. Virol.* **74**:10359–70.
- 857 32. **Nagashima S, Sasaki J, Taniguchi K.** 2003. Functional analysis of the stem-loop  
858 structures at the 5' end of the Aichi virus genome. *Virology* **313**:56–65.
- 859 33. **Piñeiro D, Martinez-Salas E.** 2012. RNA structural elements of hepatitis C virus  
860 controlling viral RNA translation and the implications for viral pathogenesis. *Viruses*  
861 **4**:2233–50.
- 862 34. **Thurner C.** 2004. Conserved RNA secondary structures in Flaviviridae genomes. *J.*  
863 *Gen. Virol.* **85**:1113–1124.
- 864 35. **Pelletier J, Sonenberg N.** 1988. Internal initiation of translation of eukaryotic mRNA  
865 directed by a sequence derived from poliovirus RNA. *Nature* **334**:320–5.
- 866 36. **Kolupaeva VG, Pestova T V, Hellen CU.** 2000. Ribosomal binding to the internal  
867 ribosomal entry site of classical swine fever virus. *RNA* **6**:1791–807.
- 868 37. **Kanamori Y, Nakashima N.** 2001. A tertiary structure model of the internal  
869 ribosome entry site (IRES) for methionine-independent initiation of translation. *RNA*  
870 **7**:266–74.

- 871 38. **Miller WA, Wang Z, Treder K.** 2007. The amazing diversity of cap-independent  
872 translation elements in the 3'-untranslated regions of plant viral RNAs. *Biochem.*  
873 *Soc. Trans.* **35**:1629–33.
- 874 39. **Simon AE, Miller WA.** 2013. 3' cap-independent translation enhancers of plant  
875 viruses. *Annu. Rev. Microbiol.* **67**:21–42.
- 876 40. **Cossons N, Faust EA, Zannis-Hadjopoulos M.** 1996. DNA polymerase delta-  
877 dependent formation of a hairpin structure at the 5' terminal palindrome of the  
878 minute virus of mice genome. *Virology* **216**:258–64.
- 879 41. **Sun Y, Chen AY, Cheng F, Guan W, Johnson FB, Qiu J.** 2009. Molecular  
880 characterization of infectious clones of the minute virus of canines reveals unique  
881 features of bocaviruses. *J. Virol.* **83**:3956–67.
- 882 42. **Perros M, Spegelaere P, Dupont F, Vanacker JM, Rommelaere J.** 1994.  
883 Cruciform structure of a DNA motif of parvovirus minute virus of mice (prototype  
884 strain) involved in the attenuation of gene expression. *J. Gen. Virol.* **75 ( Pt**  
885 **10**:2645–53.
- 886 43. **Krauskopf A, Bengal E, Aloni Y.** 1991. The block to transcription elongation at the  
887 minute virus of mice attenuator is regulated by cellular elongation factors. *Mol. Cell.*  
888 *Biol.* **11**:3515–21.
- 889 44. **Ben-Asher E, Aloni Y.** 1984. Transcription of minute virus of mice, an autonomous  
890 parvovirus, may be regulated by attenuation. *J. Virol.* **52**:266–76.
- 891 45. **Resnekov O, Aloni Y.** 1989. RNA polymerase II is capable of pausing and  
892 prematurely terminating transcription at a precise location in vivo and in vitro. *Proc.*  
893 *Natl. Acad. Sci. U. S. A.* **86**:12–6.
- 894 46. **Bohenzky RA, LeFebvre RB, Berns KI.** 1988. Sequence and symmetry  
895 requirements within the internal palindromic sequences of the adeno-associated  
896 virus terminal repeat. *Virology* **166**:316–27.
- 897 47. **Orozco BM, Hanley-Bowdoin L.** 1996. A DNA structure is required for geminivirus  
898 replication origin function. *J. Virol.* **70**:148–58.
- 899 48. **Hafner GJ, Stafford MR, Wolter LC, Harding RM, Dale JL.** 1997. Nicking and  
900 joining activity of banana bunchy top virus replication protein in vitro. *J. Gen. Virol.*  
901 **78 ( Pt 7)**:1795–9.
- 902 49. **Cheung AK.** 2006. Rolling-circle replication of an animal circovirus genome in a  
903 theta-replicating bacterial plasmid in *Escherichia coli*. *J. Virol.* **80**:8686–94.



- 904 50. **Morozov Sy, Chernov B, Merits A, Blinov V.** 1994. Computer-assisted predictions  
905 of the secondary structure in the plant virus single-stranded DNA genome. *J.*  
906 *Biomol. Struct. Dyn.* **11**:837–847.
- 907 51. **Shepherd DN, Martin DP, Varsani A, Thomson JA, Rybicki EP, Klump HH.**  
908 2006. Restoration of native folding of single-stranded DNA sequences through  
909 reverse mutations: an indication of a new epigenetic mechanism. *Arch. Biochem.*  
910 *Biophys.* **453**:108–22.
- 911 52. **Martin DP, Lefeuvre P, Varsani A, Hoareau M, Semegni J, Dijoux B, Vincent C,**  
912 **Reynaud B, Lett J.** 2011. Complex recombination patterns arising during  
913 geminivirus coinfections preserve and demarcate biologically important intra-  
914 genome interaction networks. *PLoS Pathog.* **7**:e1002203.
- 915 53. **Schultes EA, Spasic A, Mohanty U, Bartel DP.** 2005. Compact and ordered  
916 collapse of randomly generated RNA sequences. *Nat. Struct. Mol. Biol.* **12**:1130–6.
- 917 54. **Hasegawa M, Yasunaga T, Miyata T.** 1979. Secondary structure of MS2 phage  
918 RNA and bias in code word usage. *Nucleic Acids Res.* **7**:2073–9.
- 919 55. **Cardinale DJ, DeRosa K, Duffy S.** 2013. Base composition and translational  
920 selection are insufficient to explain codon usage bias in plant viruses. *Viruses*  
921 **5**:162–81.
- 922 56. **Ngandu NK, Scheffler K, Moore P, Woodman Z, Martin D, Seoighe C.** 2008.  
923 Extensive purifying selection acting on synonymous sites in HIV-1 Group M  
924 sequences. *Viol. J.* **5**:160.
- 925 57. **Cheung AK.** 2005. Detection of rampant nucleotide reversion at the origin of DNA  
926 replication of porcine circovirus type 1. *Virology* **333**:22–30.
- 927 58. **Cheng N, Mao Y, Shi Y, Tao S.** 2012. Coevolution in RNA molecules driven by  
928 selective constraints: evidence from 5S rRNA. *PLoS One* **7**:e44376.
- 929 59. **Fernández N, Fernandez-Miragall O, Ramajo J, García-Sacristán A, Bellora N,**  
930 **Eyras E, Briones C, Martínez-Salas E.** 2011. Structural basis for the biological  
931 relevance of the invariant apical stem in IRES-mediated translation. *Nucleic Acids*  
932 *Res.* **39**:8572–85.
- 933 60. **Hofacker I, Fekete M, Flamm C, Huyenen M, Rauscher S, Stolorz P, PF S.** 1998.  
934 Automatic detection of conserved RNA structure elements in complete RNA virus  
935 genomes. *Nucleic Acids Res.* **26**:3825–3836.
- 936 61. **Cheung AK.** 2004. Detection of template strand switching during initiation and  
937 termination of DNA replication of porcine circovirus. *J. Virol.* **78**:4268–77.



- 938 62. **Cheung AK.** 2004. Palindrome regeneration by template strand-switching  
939 mechanism at the origin of DNA replication of porcine circovirus via the rolling-circle  
940 melting-pot replication model. *J. Virol.* **78**:9016–29.
- 941 63. **Edgar RC.** 2004. MUSCLE: multiple sequence alignment with high accuracy and  
942 high throughput. *Nucleic Acids Res.* **32**:1792–7.
- 943 64. **Tamura K, Peterson D, Peterson N, Stecher G, Nei M, Kumar S.** 2011. MEGA5:  
944 molecular evolutionary genetics analysis using maximum likelihood, evolutionary  
945 distance, and maximum parsimony methods. *Mol. Biol. Evol.* **28**:2731–9.
- 946 65. **Wilm A, Mainz I, Steger G.** 2006. An enhanced RNA alignment benchmark for  
947 sequence alignment programs. *Algorithms Mol. Biol.* **1**:19.
- 948 66. **Semegni JY, Wamalwa M, Gaujoux R, Harkins GW, Gray A, Martin DP.** 2011.  
949 NASP: a parallel program for identifying evolutionarily conserved nucleic acid  
950 secondary structures from nucleotide sequence alignments. *Bioinformatics* **27**:2443–  
951 5.
- 952 67. **Greiner W, Neise L, Stöcker H.** 1995. Thermodynamics and statistical mechanics,  
953 p. 101. Springer-Verlag, New York.  
954 <http://www.tandar.cnea.gov.ar/~sacanell/Greiner.pdf>.
- 955 68. **Ding Y.** 2003. A statistical sampling algorithm for RNA secondary structure  
956 prediction. *Nucleic Acids Res.* **31**:7280–7301.
- 957 69. **Markham NR, Zuker M.** 2008. UNAFold: software for nucleic acid folding and  
958 hybridization., p. 3–31. *In* Keith, JM (ed.), *Methods in molecular biology* (Clifton,  
959 N.J.). Humana Press, Totowa, NJ.
- 960 70. **Tajima F.** 1989. Statistical Method for Testing the Neutral Mutation Hypothesis by  
961 DNA Polymorphism. *Genetics* **123**:585–595.
- 962 71. **Fu YX, Li WH.** 1993. Statistical tests of neutrality of mutations. *Genetics* **133**:693–  
963 709.
- 964 72. **Scheffler K, Martin DP, Seoighe C.** 2006. Robust inference of positive selection  
965 from recombining coding sequences. *Bioinformatics* **22**:2493–2499.
- 966 73. **Murrell B, Moola S, Mabona A, Weighill T, Sheward D, Kosakovsky Pond SL,  
967 Scheffler K.** 2013. FUBAR: a fast, unconstrained bayesian approximation for  
968 inferring selection. *Mol. Biol. Evol.* **30**:1196–205.
- 969 74. **Muse S V, Gaut BS.** 1994. A likelihood approach for comparing synonymous and  
970 nonsynonymous nucleotide substitution rates, with application to the chloroplast  
971 genome. *Mol. Biol. Evol.* **11**:715–24.

- 972 75. **Scheffler K, Martin DP, Seoighe C.** 2006. Robust inference of positive selection  
973 from recombining coding sequences. *Bioinformatics* **22**:2493–9.
- 974 76. **Lefeuvre P, Lett J-M, Varsani A, Martin DP.** 2009. Widely conserved  
975 recombination patterns among single-stranded DNA viruses. *J. Virol.* **83**:2697–707.
- 976 77. **Julian L, Piasecki T, Chrzastek K, Walters M, Muhire B, Harkins GW, Martin DP,**  
977 **Varsani A.** 2013. Extensive recombination detected among beak and feather  
978 disease virus isolates from breeding facilities in Poland. *J. Gen. Virol.* **94**:1086–95.
- 979 78. **Padidam M, Sawyer S, Fauquet CM.** 1999. Possible emergence of new  
980 geminiviruses by frequent recombination. *Virology* **265**:218–25.
- 981 79. **Shackelton L a, Hoelzer K, Parrish CR, Holmes EC.** 2007. Comparative analysis  
982 reveals frequent recombination in the parvoviruses. *J. Gen. Virol.* **88**:3294–301.
- 983 80. **Navas-Castillo J, Sánchez-Campos S, Noris E, Louro D, Accotto GP, Moriones**  
984 **E.** 2000. Natural recombination between Tomato yellow leaf curl virus-is and Tomato  
985 leaf curl virus. *J. Gen. Virol.* **81**:2797–801.
- 986 81. **Kosakovsky Pond SL, Posada D, Gravenor MB, Woelk CH, Frost SDW.** 2006.  
987 GARD: a genetic algorithm for recombination detection. *Bioinformatics* **22**:3096–8.
- 988 82. **Pond SLK, Frost SDW, Muse S V.** 2005. HyPhy: hypothesis testing using  
989 phylogenies. *Bioinformatics* **21**:676–9.
- 990 83. **Hasegawa M, Kishino H, Yano T.** 1985. Dating of the human-ape splitting by a  
991 molecular clock of mitochondrial DNA. *J. Mol. Evol.* **22**:169–174.
- 992 84. **Muse S V.** 1995. Evolutionary analyses of DNA sequences subject to constraints of  
993 secondary structure. *Genetics* **139**:1429–39.
- 994 85. **Martin DP, Lemey P, Lott M, Moulton V, Posada D, Lefeuvre P.** 2010. RDP3: a  
995 flexible and fast computer program for analyzing recombination. *Bioinformatics*  
996 **26**:2462–3.
- 997 86. **Guindon S, Dufayard J-F, Lefort V, Anisimova M, Hordijk W, Gascuel O.** 2010.  
998 New algorithms and methods to estimate maximum-likelihood phylogenies:  
999 assessing the performance of PhyML 3.0. *Syst. Biol.* **59**:307–21.
- 1000 87. **McCormack JC, Simon AE.** 2004. Biased hypermutagenesis associated with  
1001 mutations in an untranslated hairpin of an RNA virus. *J. Virol.* **78**:7813–7.
- 1002 88. **Staplin WR, Miller WA.** 2008. In vivo analyses of viral RNA translation. *Methods*  
1003 *Mol. Biol.* **451**:99–112.

- 1004 89. **Simmonds P, Smith DB.** 1999. Structural constraints on RNA virus evolution. *J.*  
1005 *Virol.* **73**:5787–94.
- 1006 90. **Frederico L, Kunkel T, Shaw B.** 1990. A sensitive genetic assay for the detection  
1007 of cytosine deamination: determination of rate constants and the activation energy.  
1008 *Biochemistry* **29**:2532–2537.
- 1009 91. **Xia X, Yuen KY.** 2005. Differential selection and mutation between dsDNA and  
1010 ssDNA phages shape the evolution of their genomic AT percentage. *BMC Genet.*  
1011 **6**:20.
- 1012 92. **Monjane AL, Pande D, Lakay F, Shepherd DN, van der Walt E, Lefeuvre P, Lett**  
1013 **J-M, Varsani A, Rybicki EP, Martin DP.** 2012. Adaptive evolution by recombination  
1014 is not associated with increased mutation rates in Maize streak virus. *BMC Evol.*  
1015 *Biol.* **12**:252.
- 1016 93. **Fisher RA.** 1922. On the Interpretation of  $\chi^2$  from Contingency Tables, and the  
1017 Calculation of P. *J. R. Stat. Soc.* **85**:87.
- 1018 94. **Golden M, Martin D.** 2013. DOOSS: a tool for visual analysis of data overlaid on  
1019 secondary structures. *Bioinformatics* **29**:271–2.
- 1020 95. **Ray SS, Pal SK.** 2013. RNA secondary structure prediction using soft computing.  
1021 *IEEE/ACM Trans. Comput. Biol. Bioinform.* **10**:2–17.
- 1022 96. **Willwand K, Hirt B.** 1991. The minute virus of mice capsid specifically recognizes  
1023 the 3' hairpin structure of the viral replicative-form DNA: mapping of the binding site  
1024 by hydroxyl radical footprinting. *J. Virol.* **65**:4629–35.
- 1025 97. **Wright E, Heckel T, Groenendijk J, Davies J, Boulton M.** 1997. Splicing features  
1026 in maize streak virus virion- and complementary-sense gene expression. *Plant J.*  
1027 **12**:1285–1297.
- 1028 98. **Warf MB, Berglund JA.** 2010. Role of RNA structure in regulating pre-mRNA  
1029 splicing. *Trends Biochem. Sci.* **35**:169–78.
- 1030 99. **Munroe SH.** 1984. Secondary structure of splice sites in adenovirus mRNA  
1031 precursors. *Nucleic Acids Res.* **12**:8437–56.
- 1032 100. **De Villiers E-M, Borkosky SS, Kimmel R, Gunst K, Fei J-W.** 2011. The diversity  
1033 of torque teno viruses: in vitro replication leads to the formation of additional  
1034 replication-competent subviral molecules. *J. Virol.* **85**:7284–95.

1035

1036 **Table and Figure legends**

1037 **Table 1. List of the 23 large datasets analysed.**

1038

1039 **Table 2. Tajima's D and Fu and Li F statistics for paired and unpaired genomic site**  
1040 **alignments.**

1041

1042 **Table 3. Comparison of synonymous substitution rates at paired and unpaired**  
1043 **codon-sites.**

1044

1045 **Table 4. Tests for associations between paired-sites and complementarily**  
1046 **coevolving sites.**

1047

1048 **Table 5. Summary of the various tests carried on the 23 different ssDNA virus**  
1049 **genome sequence datasets.**

1050

1051 **Fig. 1. Secondary-structure map of plant-infecting ssDNA viruses.**

1052 Genome organisation maps of plant-infecting ssDNA viruses. In each map, the arcs  
1053 represent the coordinates of identified structural elements within "high confidence structure  
1054 sets" (HCSS). These highly conserved structural elements are those primarily responsible  
1055 for the estimated structural stability of the analysed genomes being greater than that of  
1056 randomised sequences with identical nucleotide compositions. The ten structures  
1057 collectively displaying the greatest degrees of base-pairing conservation, lowest  
1058 associated synonymous substitution rates and greatest degrees of complementary  
1059 coevolution between paired nucleotides are shown using arcs in cyan and magenta (to  
1060 distinguish the two complementary parts of the stem sequences). All remaining structures  
1061 are shown using black arcs. Black arrows indicate examples of currently uncharacterised  
1062 but likely biologically functional structures that are apparently conserved across multiple

1063 datasets (coloured in green and brown when these were not ranked among the top ten  
1064 within their respective HCSS). The known secondary-structural elements at the virion  
1065 strand origins of replication are indicated by a star symbol at the 12 o'clock position of the  
1066 genomes. Numbers in brackets indicate the lengths of the genomes in kilobases (kb) and  
1067 the ratio of the numbers of high confidence structures over the total numbers of predicted  
1068 secondary-structures. Italicized abbreviations of gene names: *rep*=replication associated  
1069 protein; *cp* = coat protein; *mp* = movement protein; *clink* = cell cycle link protein; *nsp* =  
1070 nuclear shuttle protein; *ren* = replication enhancer protein; *trap* = transcription activator  
1071 protein.

1072

1073 **Fig. 2. Secondary-structure map of animal-infecting ssDNA viruses.**

1074 Genome organisation maps of animal-infecting ssDNA viruses. In each map, the arcs (for  
1075 circular genomes) and vertical lines (for linear genomes) represent the coordinates of  
1076 identified structural elements within "high confidence structure sets" (HCSS). The ten  
1077 structures collectively displaying the greatest degrees of base-pairing conservation, lowest  
1078 associated synonymous substitution rates and greatest degrees of complementary  
1079 coevolution between paired nucleotides are shown using arcs in cyan and magenta (to  
1080 distinguish the two complementary parts of the stem sequences). All remaining structures  
1081 are shown using black arcs/vertical lines. Black arrows indicate examples of currently  
1082 uncharacterised but likely biologically functional structures that are apparently conserved  
1083 across multiple datasets (coloured in green and brown when these were not ranked  
1084 among the top ten structures within their respective HCSS). The known secondary-  
1085 structural elements at the virion strand origins of replication are indicated by a star symbol

1086 at the 12 o'clock position of the genomes. Numbers in brackets indicate the lengths of the  
1087 genomes in kilobases (kb) and the ratio of the numbers of high confidence structures over  
1088 the total numbers of predicted secondary-structures. Italicized abbreviations represent the  
1089 gene names encoding the following proteins: *rep*=replication associated protein; *cp* = coat  
1090 protein; *ns1* = large non-structural protein; *np1* = small non-structural protein; *vp1* = major  
1091 virion/viral protein; *vp2* = minor virion/viral protein and *ORF* = unnamed open reading  
1092 frame.

1093

1094 **Fig. 3. Secondary-structure associated with the intron of the mastrevirus movement**  
1095 **protein gene.**

1096 A secondary-structure associated with the movement protein gene intron was predicted in  
1097 all five mastrevirus datasets. This structure is highly conserved and contains splice donor  
1098 and acceptor sites (indicated by arrows), as well as, in the case of the GeminiMSV,  
1099 GeminiPanSV and GeminiWDV, a likely lariat sequences (outlined in pink). The similarities  
1100 between these structures include homologous stem-loop structures conserved in all but  
1101 GeminiWDV (highlighted in blue and yellow), a highly conserved stem-structure found in  
1102 both GeminiMSV and GeminiPanSV, and conserved sequences in the stems of  
1103 GeminiTYDV-CpCV and GeminiCpCDV. The rank ratio shows the actual rank of a  
1104 structure over the total number of structures predicted in the high confidence structure set  
1105 (HCSS). This structure is highly ranked in GeminiCpCDV and GeminiTYDV-CpCV (ranked  
1106 7<sup>th</sup> out of 25 structures in HCSS and 10<sup>th</sup> out of 47 structures in the HCSS set  
1107 respectively). In case of GeminiCpCDV, base-pairing interactions displaying significant  
1108 associated complementary coevolution ( $p$ -value <0.05) are represented by a red line

1109 where the degree of redness reflects the  $p$ -value. Whereas nucleotide sequence variability  
 1110 is reflected by a sequence logo at each position, each position is also associated with a  
 1111 colour ranging from blue to green depicting the rate of synonymous substitutions of the  
 1112 codon site at which the nucleotide is located. Low synonymous substitution rates are  
 1113 observable in the stem region in all datasets, indicating that there is a high degree of  
 1114 conservation at these particular sites. Although the sequence of this structure is divergent  
 1115 in all five mastrevirus datasets, it is plausible that this structure has some function during  
 1116 splicing of the movement protein intron.

1117

1118 **Fig. 4. Secondary-structure associated with the 3' end of the begomovirus coat**  
 1119 **protein gene.**

1120 A secondary-structure with a potential role in transcriptional termination was predicted at  
 1121 the end of the coat protein gene of the begomovirus datasets, GeminiTYLCV,  
 1122 GeminiEACMV and GeminiMYVYV. In all these, the structure has a stop codon and a  
 1123 stem-loop containing a polyadenylation signal (the complementary polyadenylation  
 1124 signalling sequences within the stem-loops are in bold text). A common stem-loop  
 1125 structure between the GeminiEACMV and GeminiMYVYV dataset is highlighted in yellow.  
 1126 Nucleotide logos and colours respectively indicate degrees of sequence variability and  
 1127 associated synonymous nucleotide substitution rates as outlined in Fig. 3. Nucleotides  
 1128 falling outside genes are shaded grey. Base-pairing interactions displaying significant  
 1129 associated complementary coevolution ( $p$ -value  $<0.05$ ) are represented by a red line  
 1130 where the degree of redness reflects the  $p$ -value. The rank ratio shows the actual rank of a  
 1131 structure over the total number of structures predicted in the high confidence structure set.



1132

1133 **Fig. 5. Parvovirus secondary-structures predicted at the start of genes.**

1134 Secondary-structures predicted at the start of genes represented in the parvovirus  
1135 datasets ParvoAAV, ParvoHBoV and ParvoMPV are shown. These include those  
1136 spanning the start of the large non-structural proteins (*ns1*; P1), the major viral/virion  
1137 proteins (*vp1*; P2), the small non-structural protein (*np1*; P3) and the minor viral/virion  
1138 proteins (*vp2*, P4). Nucleotide logos and colours respectively indicate degrees of  
1139 sequence variability and associated synonymous nucleotide substitution rates as outlined  
1140 in Fig. 3. Base-pairing interactions displaying significant associated complementary  
1141 coevolution ( $p$ -value  $< 0.05$ ) are represented by a line where the degree of redness reflects  
1142 the  $p$ -value. The rank ratio shows the actual rank of a structure over the total number of  
1143 structures predicted in the high confidence structure set. The ParvoMPV IR-*ns1* stem-loop  
1144 involved in transcription attenuation is highlighted in grey. In the depicted structures start  
1145 codons are consistently located either within or immediately adjacent to an unpaired loop  
1146 or bulge, which might enhance the accessibility of these codons during transcription or  
1147 translation.

1148

1149 **Fig. 6. Conserved circovirus stem-loop structure within the intergenic region.**

1150 A stem-loop structures that are highly conserved within the intergenic regions of the two  
1151 circovirus datasets, CircoCoCV and CircoBFDV, is shown. Nucleotide logos reflect  
1152 degrees of sequence variability. The stems of these structures have high GC-contents and  
1153 display clear evidence of complementary coevolution between base-paired nucleotides  
1154 within the CircoBFDV dataset (base-pairing interactions displaying significant associated



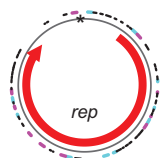
1155 complementary coevolution with  $p$ -value  $<0.05$  are represented by red lines where the  
1156 degree of redness reflects the  $p$ -value). Additionally, a pentanucleotide loop (highlighted in  
1157 green) is highly conserved at the top of the stem in both datasets.

1158

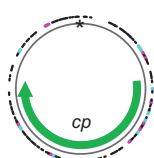
1159 **Fig. 7. Anellovirus highly conserved intergenic T-shaped structures.**

1160 A T-shaped structure predicted within the intergenic region (IR) of two divergent  
1161 anellovirus datasets, AnelloTTSuV1 and AnelloTTSuV2 is shown. These structures have  
1162 homologous 17-nucleotide long sequences on top of the “T” (highlighted in sky blue) and  
1163 similar sequences in the stem (highlighted in green, yellow and blue). The homologue to  
1164 these structures in the even more divergent AnelloTTV dataset has a stem-loop rather  
1165 than a “T” configuration. It shares similar sequences (highlighted using yellow and blue)  
1166 with the ones found in the other anellovirus datasets. In the AnelloTTV structure, base-  
1167 pairing interactions displaying significant associated complementary coevolution ( $p$ -value  
1168  $<0.05$ ) are represented by a red line where the degree of redness reflects the  $p$ -value.  
1169 Nucleotide logos reflect the degree of sequence diversity at individual sites.

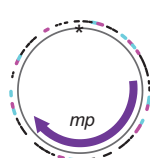
## Nanovirus



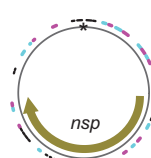
**NanoBBTVR**  
[EF546807]  
(1.1 kb, 31/76)



**NanoBBTVS**  
[AB113661]  
(1.1 kb, 49/80)



**NanoBBTVM**  
[JF957665]  
(1.1 kb, 32/71)

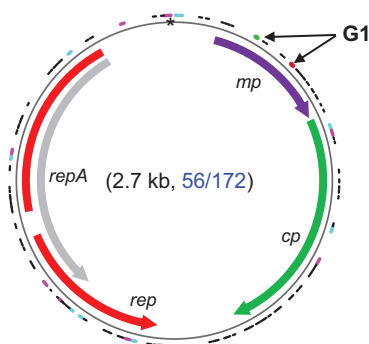


**NanoBBTVN**  
[AF238879]  
(1.1 kb, 19/83)

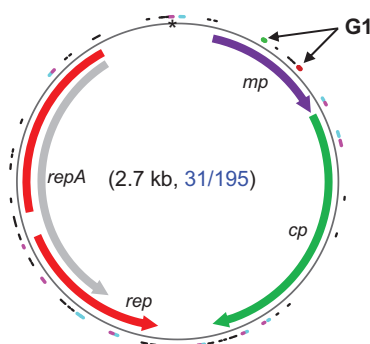


**NanoBBTVC**  
[EF520722]  
(1 kb, 24/69)

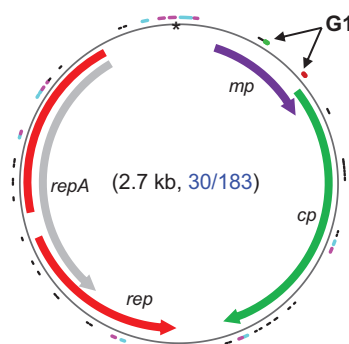
## Geminivirus



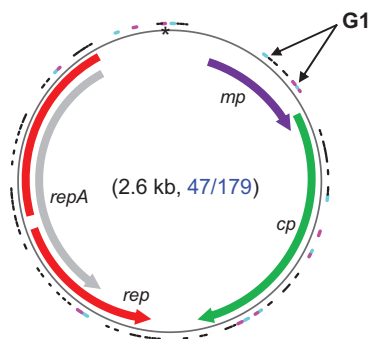
**GeminiMSV** [Y00514]



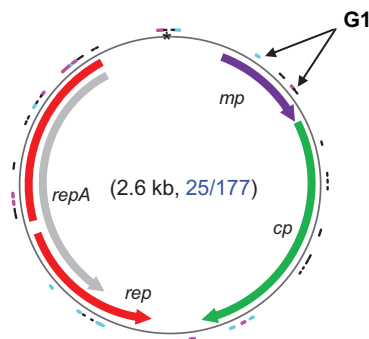
**GeminiPanSV** [L39638]



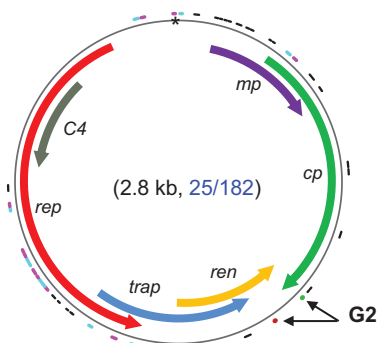
**GeminiWDV** [FJ620684]



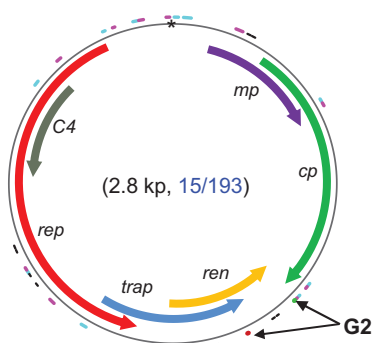
**GeminiTYDV-CpCV** [JN989413]



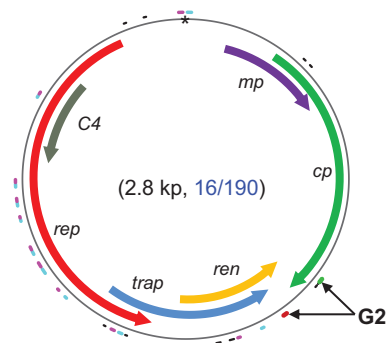
**GeminiCpCDV** [AM849097]



**GeminiTYLCV** [EU847740]

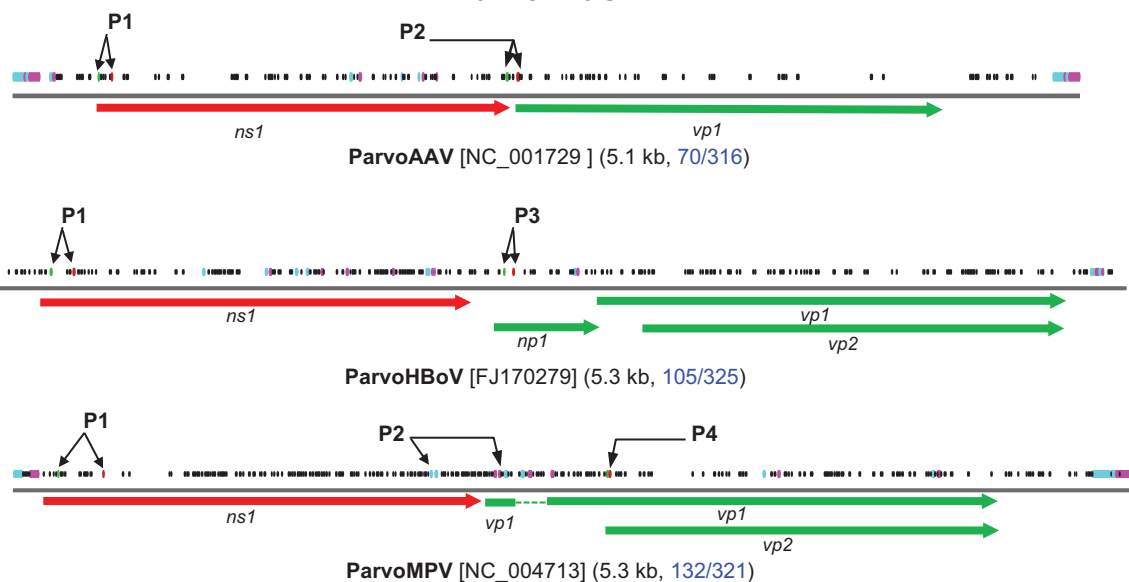


**GeminiEACMV** [AJ717558]

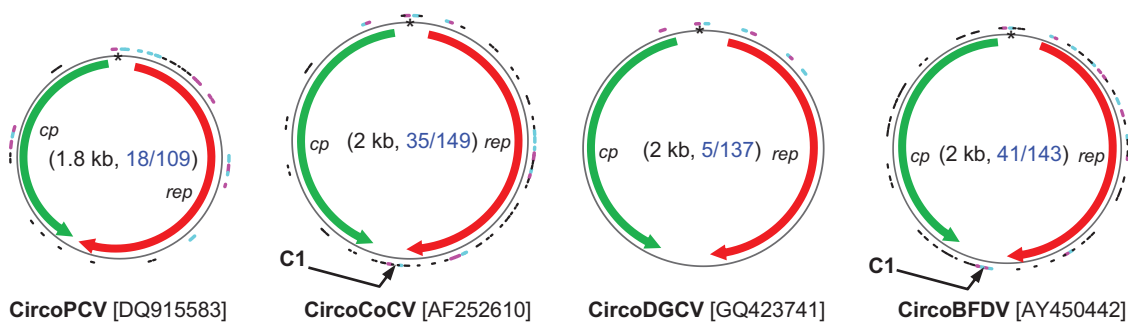


**GeminiMYVVYV** [EF185318]

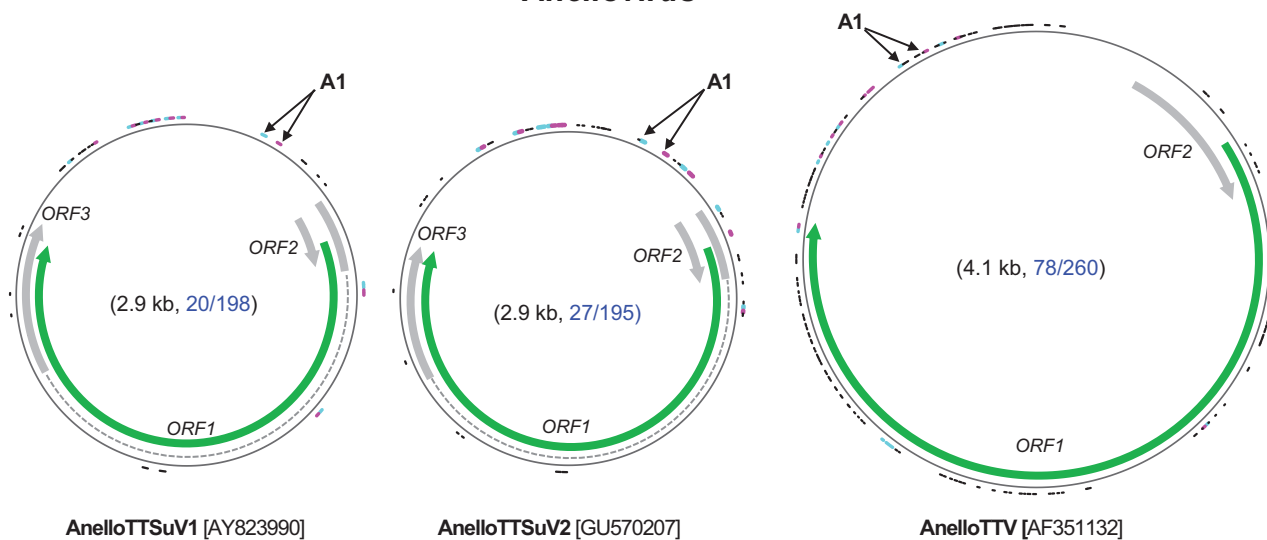
## Parvovirus

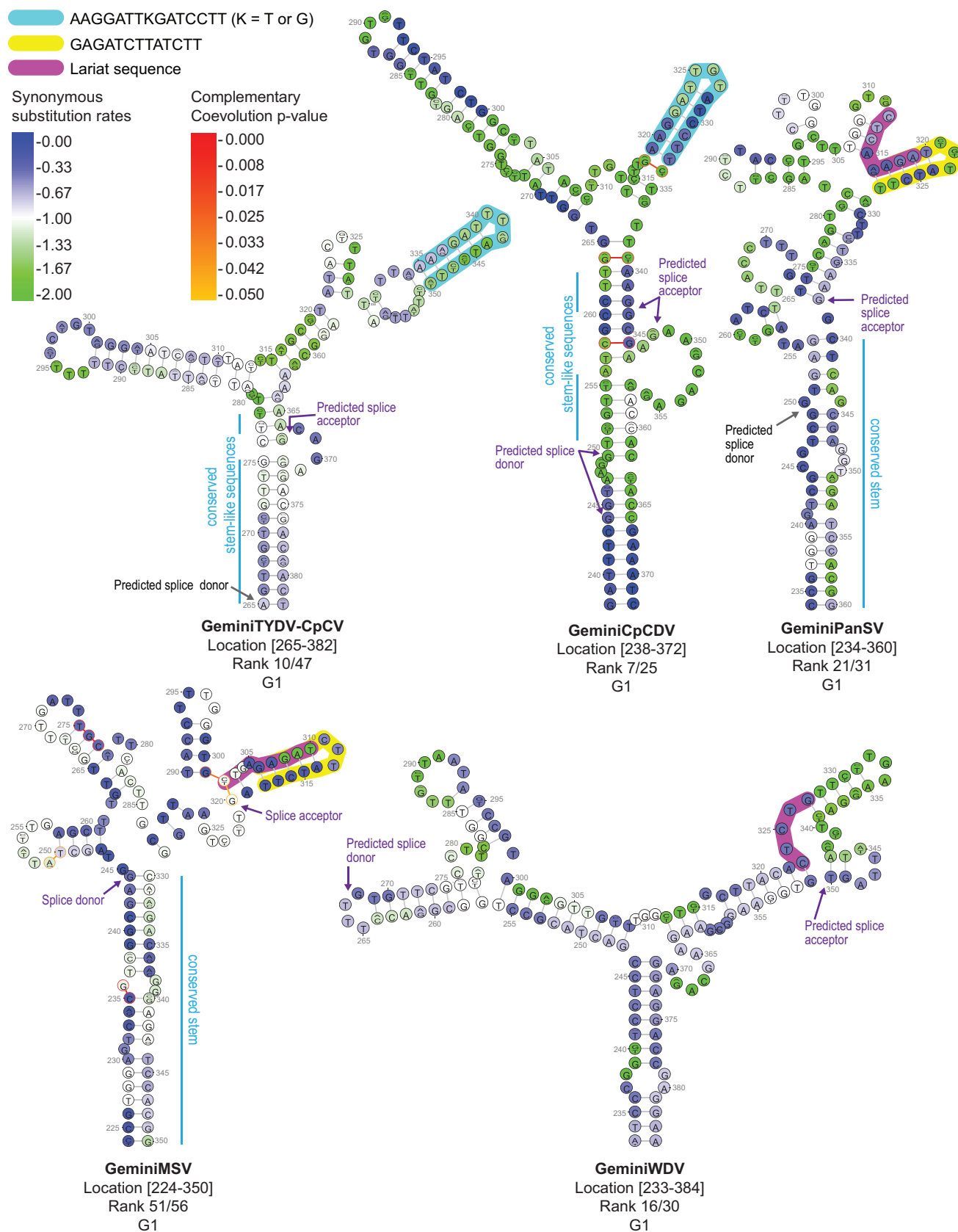


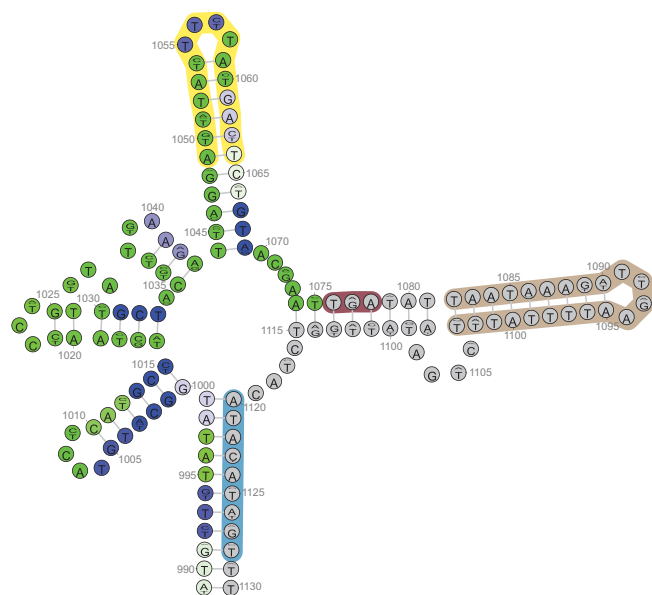
## Circovirus



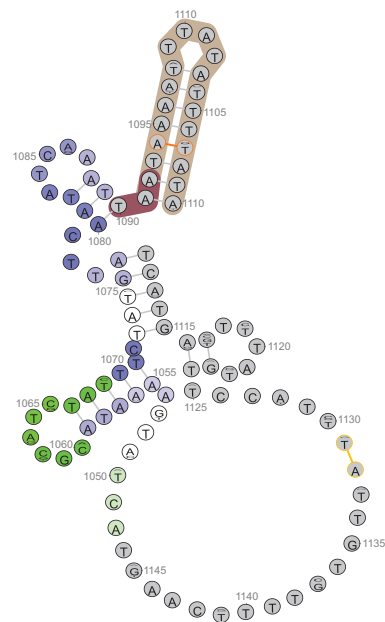
## Anellovirus



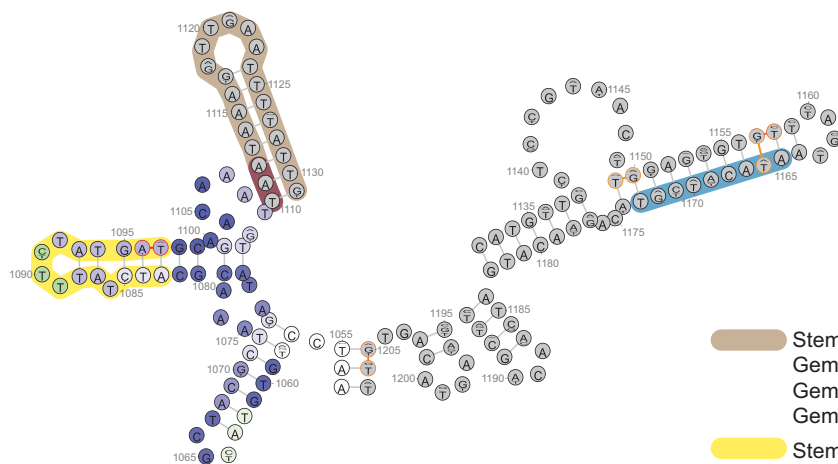




**GeminiMYVVY**  
Location [989-1130]  
Rank 13/16  
G2



**GeminiTYLCV**  
Location [1048-1146]  
rank 21/25  
G2

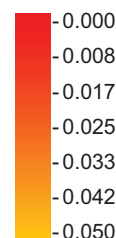


**GeminiEACMV**  
Location [1053-1207]  
Rank 14/15  
G2

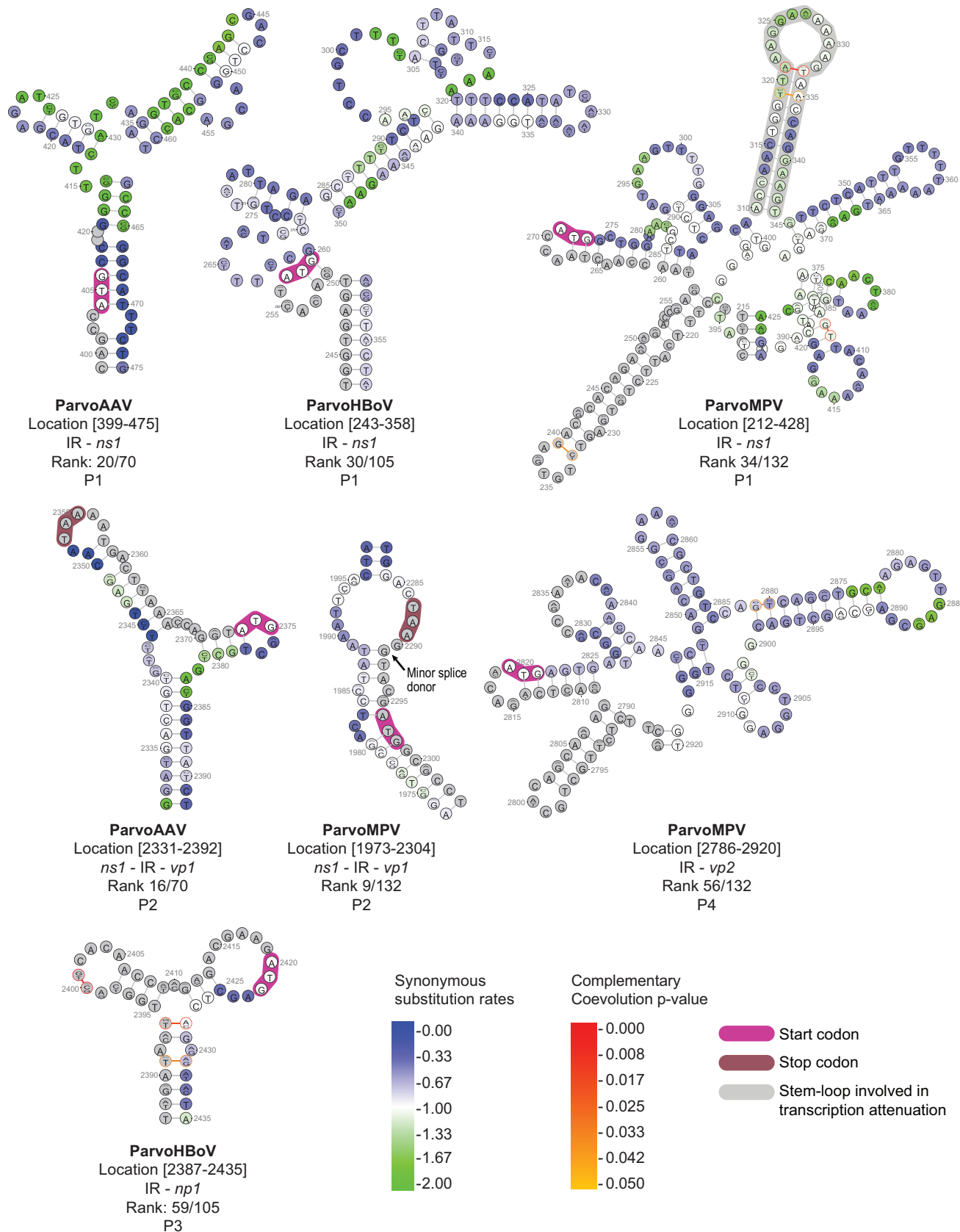
Synonymous  
substitution rates

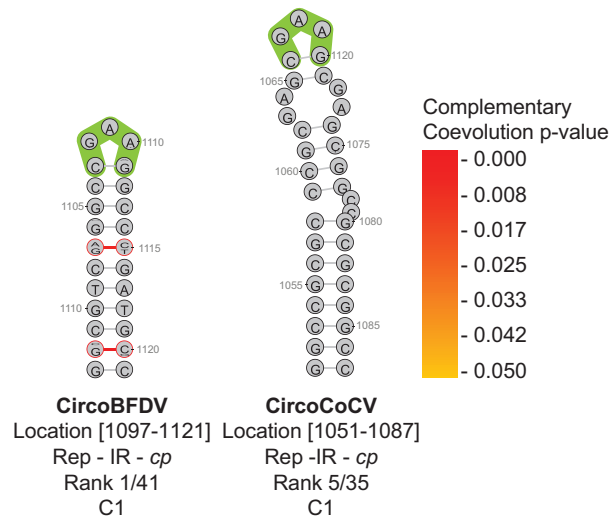


Complementary  
Coevolution p-value

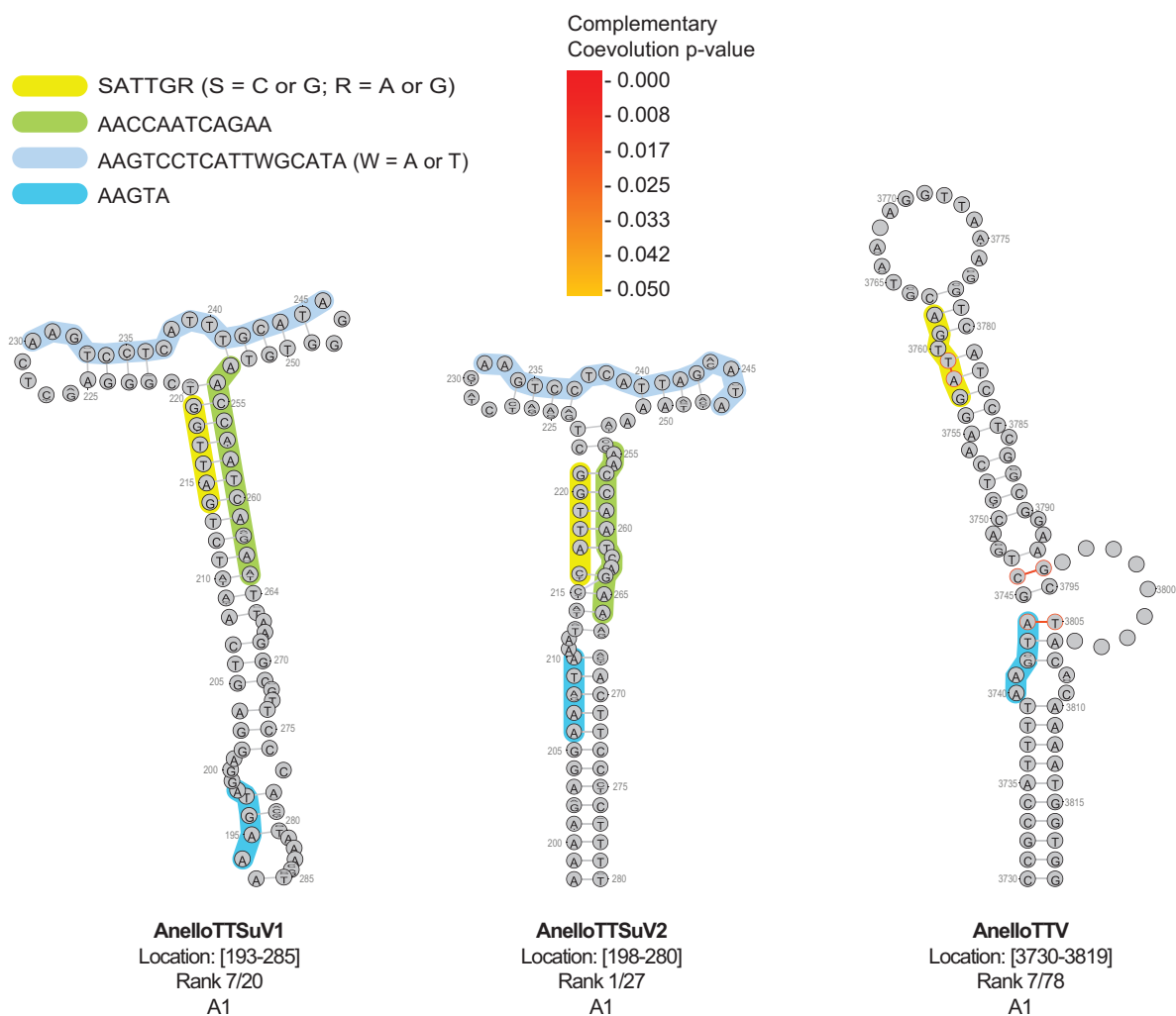


- Stem loop  
GeminiMYVVY **TAATAAAGATTGAATTTTATTT**  
GeminiTYLCV **TAATAAAATTTATATTTTAT**  
GeminiEACMV **TAATAAAGGTTGAATTTTATTG**
- Stem loop  
ATYTATTTCTATGAT (Y= C or T)
- ATACAYTGT (Y=C or T)
- Stop codon











**Table 1. List of the 23 of large datasets obtained.**

	<b>Name<sup>a</sup></b>	<b>Families</b>	<b>Constituent virus species</b>	<b>Size<sup>b</sup></b>
1	CircoPCV	Circoviridae	<i>Porcine circovirus 2</i>	519
2	CircoCoCV		<i>Columbid circovirus</i>	36
3	CircoDGCV		<i>Duck circovirus, Goose circovirus, Muscovy duck circovirus, Cygnus olor circovirus</i>	49
4	CircoBFDV		<i>Beak and feather disease virus</i>	184
5	AnelloTTSuV1	Anelloviridae	<i>Torque teno sus virus 1</i>	21
6	AnelloTTSuV2		<i>Torque teno sus virus 2, Porcine torque teno virus 2</i>	44
7	AnelloTTV	Parvoviridae	<i>Torque teno virus</i>	22
8	ParvoAAV		<i>Adeno-associated virus, Human bocavirus - 2, 3,4; Porcine bocavirus 1,2; Gorilla bocavirus, Bovine parvovirus 1, Canine minute virus</i>	34
9	ParvoHBoV		<i>Mouse parvovirus 4, Rat minute virus, Mouse parvovirus, Minute virus, Lull virus, Hamster parvovirus</i>	21
10	ParvoMPV		<i>Mouse parvovirus, Minute virus, Lull virus, Hamster parvovirus</i>	26
11	NanoBBTV-R	Nanoviridae	<i>Banana bunchy top virus component R</i>	221
12	NanoBBTV-S		<i>Banana bunchy top virus component S</i>	189
13	NanoBBTV-M		<i>Banana bunchy top virus component M</i>	150
14	NanoBBTV-N		<i>Banana bunchy top virus component N</i>	148
15	NanoBBTV-C	Geminiviridae	<i>Banana bunchy top virus component C</i>	122
16	GeminiMSV		<i>Maize streak virus</i>	759
17	GeminiWDV		<i>Wheat dwarf virus</i>	138
18	GeminiPanSV		<i>Panicum streak virus</i>	41
19	GeminiTYDV-CpCV		<i>Tobacco yellow dwarf virus, Chickpea chlorosis virus, Chickpea yellows virus</i>	41
20	GeminiCpCDV		<i>Chickpea chlorotic dwarf virus</i>	43
21	GeminiTYLCV		<i>Tomato yellow leaf curl virus</i>	228
22	GeminiEACMV		<i>East African cassava mosaic virus, South African cassava mosaic virus</i>	146
23	GeminiMYVVV		<i>Malvastrum yellow vein Yunnan virus, Cotton leaf curl Multan virus isolate, Bhendi yellow vein India virus</i>	254

<sup>a</sup> The name of the dataset is made up of the prefix of its family and the abbreviation of the main virus species it contains.

<sup>b</sup> The number of full genome sequences in the dataset

**Table 2. Tajima's D – Fu and Li statistics for paired and unpaired genomic site alignments**

	Large dataset	Tajima's D			Fu & Li's F		
		Paired <sup>a</sup>	Permuted unpaired <sup>b</sup>	p-value	Paired <sup>c</sup>	Permuted unpaired <sup>d</sup>	p-value
<b>1</b>	CircoPCV	-2.08	-1.90	0.06	-5.33	-4.91	0.13
<b>2</b>	CircoCoCV	-2.48	-1.75	< 0.01	-4.98	-3.17	< 0.01
<b>3</b>	CircoDGCV	1.44	0.71	0.99	1.40	0.58	0.96
<b>4</b>	CircoBFDV	-1.53	-1.33	< 0.01	-2.54	-1.77	< 0.01
<b>5</b>	AnelloTTSuV1	-0.72	-1.09	1	-0.19	-1.09	1
<b>6</b>	AnelloTTSuV2	-0.96	-0.95	0.54	-1.09	-1.38	0.92
<b>7</b>	AnelloTTV	-0.10	-0.09	0.5	-0.95	-0.95	0.49
<b>8</b>	ParvoAAV	-0.56	-0.55	0.53	0.80	0.67	0.85
<b>9</b>	ParvoHBoV	-1.09	-1.02	0.07	0.25	0.18	0.78
<b>10</b>	ParvoMPV	-0.22	-0.12	0.08	-0.31	-0.01	0.01
<b>11</b>	NanoBBTVR	-1.31	-1.21	0.28	-3.08	-2.27	0.03
<b>12</b>	NanoBBTVS	-0.96	-0.65	< 0.01	-2.79	-2.00	0.02
<b>13</b>	NanoBBTVM	-0.77	-0.38	0.05	-2.32	-1.74	0.14
<b>14</b>	NanoBBTVN	-1.61	-1.45	0.13	-4.26	-3.36	0.02
<b>15</b>	NanoBBTVC	-1.44	-0.80	< 0.01	-5.03	-3.28	< 0.01
<b>16</b>	GeminiMSV	-2.02	-1.72	< 0.01	-5.34	-3.26	< 0.01
<b>17</b>	GeminiWDV	-1.26	-0.61	< 0.01	-4.03	-2.99	0.04
<b>18</b>	GeminiPanSV	-0.91	-0.61	< 0.01	-0.53	-0.05	< 0.01
<b>19</b>	GeminiTYDV-CpCV	-1.28	-0.55	< 0.01	-2.30	-0.29	< 0.01
<b>20</b>	GeminiCpCDV	-0.88	-0.10	< 0.01	-0.71	0.19	< 0.01
<b>21</b>	GeminiTYLCV	-1.67	-1.71	0.61	-3.44	-3.23	0.26
<b>22</b>	GeminiEACMV	-1.33	-0.83	< 0.01	-2.58	-1.31	< 0.01
<b>23</b>	GeminiMYVYV	-1.23	-0.78	< 0.01	-3.25	-1.28	< 0.01

<sup>a</sup>Tajima's D for paired-sites alignments corresponding to the HCSS.

<sup>b</sup>Average Tajima's D for 100 permuted alignments sampled for the unpaired-sites.

<sup>c</sup>Fu and Li's F for paired-sites alignments corresponding to the HCSS.

<sup>d</sup>Average Fu and Li's F for 100 permuted alignments sampled for the unpaired-sites

**Table 3. Comparison of synonymous substitution rates at paired and unpaired codon-sites.**

	Datasets	Genes studied	Number of sequences	PARRIS <sup>a</sup>	FUBAR <sup>b</sup>
1	CircoPCV	<i>rep, cp</i>	30, 29	<i>rep</i>	<i>rep</i>
2	CircoCoCV	<i>rep, cp</i>	30, 30	<i>rep</i>	<i>rep</i>
3	CircoDGCV	<i>rep, cp</i>	30, 30	<i>rep</i>	<i>rep</i>
4	CircoBFDV	<i>rep, cp</i>	30, 29	<i>rep, cp</i>	<i>rep, cp</i>
5	AnelloTTSuV1	ORF1	17	ORF1	ORF1
6	AnelloTTSuV2	ORF1	30	ORF1	ORF1
7	AnelloTTV	ORF1	21	ORF1	ORF1
8	ParvoAAV	<i>ns1, vp1</i>	23, 30	<i>ns1, vp1</i>	<i>ns1, vp1</i>
9	ParvoHBoV	<i>ns1, np1</i>	21, 21	<i>ns1</i>	<i>ns1</i>
10	ParvoMPV	<i>ns1, vp2</i>	25, 18	<i>ns1</i>	<i>ns1, vp2</i>
11	NanoBBTV-R	<i>rep</i>	28	<i>rep</i>	<i>rep</i>
12	NanoBBTV-S	<i>cp</i>	29	<i>cp</i>	<i>cp</i>
13	NanoBBTV-M	<i>mp</i>	27	<i>mp</i>	<i>mp</i>
14	NanoBBTV-N	<i>nsp</i>	27	-	-
15	NanoBBTV-C	<i>clink</i>	30	-	-
16	GeminiMSV	<i>rep, cp, mp</i>	30, 30, 30	<i>rep, cp</i>	<i>rep, cp, mp</i>
17	GeminiWDV	<i>rep, cp, mp</i>	30, 30, 30	<i>cp, mp</i>	<i>cp, mp</i>
18	GeminiPanSV	<i>rep, cp, mp</i>	30, 30, 30	-	-
19	GeminiTYDV-CpCV	<i>rep, cp, mp</i>	30, 30, 30	<i>cp</i>	<i>cp</i>
20	GeminiCpCDV	<i>rep, cp, mp</i>	30, 30, 30	<i>rep, cp, mp</i>	<i>rep, cp, mp</i>
21	GeminiTYLCV	<i>rep, cp</i>	27, 30	-	-
22	GeminiEACMV	<i>rep, cp</i>	30, 30	<i>rep</i>	<i>rep</i>
23	GeminiMYVYV	<i>rep, cp</i>	29, 28	<i>rep</i>	<i>rep</i>

<sup>a</sup>Gene alignments in which the synonymous substitution rates (computed using PARRIS) at paired codon sites are significantly (Mann Whitney U test p-value < 0.05) lower than those at unpaired codon sites.

<sup>b</sup>Gene alignments in which the synonymous substitution rates (computed using FUBAR) at paired codon sites are significantly (Mann Whitney U test p-value < 0.05) lower than those at unpaired codon sites.

**Table 4. Association between paired sites and complementarily coevolving sites.**

	<b>Dataset</b>	<b>Chi-squared value</b>	<b>p-values</b>
<b>1</b>	CircoPCV	190.9307	$4.20 \times 10^{-14}$
<b>2</b>	CircoCoCV	0.2272	0.14
<b>3</b>	CircoDGCV	143.2324	$3.15 \times 10^{-14}$
<b>4</b>	CircoBFDV	62.5998	$1.59 \times 10^{-13}$
<b>5</b>	ParvoAAV	185.5472	$4.08 \times 10^{-14}$
<b>6</b>	ParvoHBoV	96.656	$2.13 \times 10^{-14}$
<b>7</b>	ParvoMPV	137.077	$3.02 \times 10^{-14}$
<b>8</b>	AnelloTTSuV1	117.9971	$2.60 \times 10^{-14}$
<b>9</b>	AnelloTTSuV2	38.2243	$2.41 \times 10^{-08}$
<b>10</b>	AnelloTTV	70.6212	$1.55 \times 10^{-14}$
<b>11</b>	NanoBBTV-R	107.8986	$2.37 \times 10^{-14}$
<b>12</b>	NanoBBTV-S	20.398	$1.28 \times 10^{-04}$
<b>13</b>	NanoBBTV-M	49.9491	$7.88 \times 10^{-11}$
<b>14</b>	NanoBBTV-N	48.2752	$1.79 \times 10^{-10}$
<b>15</b>	NanoBBTV-C	21.1911	$8.81 \times 10^{-05}$
<b>16</b>	GeminiMSV	212.2187	$4.67 \times 10^{-14}$
<b>17</b>	GeminiWDV	89.9702	$1.98 \times 10^{-14}$
<b>18</b>	GeminiPanSV	82.3437	$1.81 \times 10^{-14}$
<b>19</b>	GeminiTYDV-CpCV	28.6975	$2.43 \times 10^{-06}$
<b>20</b>	GeminiCpCDV	98.1122	$2.16 \times 10^{-14}$
<b>21</b>	GeminiTYLCV	159.2665	$3.50 \times 10^{-14}$
<b>22</b>	GeminiEACMV	175.6639	$3.86 \times 10^{-14}$
<b>23</b>	GeminiMYVYV	364.9167	$8.03 \times 10^{-14}$

**Table 5. Summary of results.**

	<b>Dataset</b>	<b>&gt;5 NASP structures<sup>a</sup></b>	<b>dS paired codon sites&lt; dS unpaired codon sites<sup>b</sup></b>	<b>Selection at paired sites<sup>c</sup></b>	<b>Complementary coevolution<sup>d</sup></b>
<b>1</b>	CircoPCV	+	+	-	+
<b>2</b>	CircoCoCV	+	+	+	-
<b>3</b>	CircoDGCV	-	+	-	+
<b>4</b>	CircoBFDV	+	+	+	+
<b>5</b>	AnelloTTSuV1	+	+	-	+
<b>6</b>	AnelloTTSuV2	+	+	-	+
<b>7</b>	AnelloTTV	+	+	-	+
<b>8</b>	ParvoAAV	+	+	-	+
<b>9</b>	ParvoHBoV	+	+	-	+
<b>10</b>	ParvoMPV	+	+	+	+
<b>11</b>	NanoBBTV-R	+	+	+	+
<b>12</b>	NanoBBTV-S	+	+	+	+
<b>13</b>	NanoBBTV-M	+	+	+	+
<b>14</b>	NanoBBTV-N	+	-	+	+
<b>15</b>	NanoBBTV-C	+	-	+	+
<b>16</b>	GeminiMSV	+	+	+	+
<b>17</b>	GeminiWDV	+	+	+	+
<b>18</b>	GeminiPanSV	+	-	+	+
<b>19</b>	GeminiTYDV-CpCV	+	+	+	+
<b>20</b>	GeminiCpCDV	+	+	+	+
<b>21</b>	GeminiTYLCV	+	-	-	+
<b>22</b>	GeminiEACMV	+	+	+	+
<b>23</b>	GeminiMYVYV	+	+	+	+

<sup>a</sup> datasets that had more than 5 structures significantly conserved in all lineages are given a "+" sign

<sup>b</sup> datasets in which at least for one gene alignment the synonymous substitution rates at paired codon sites were significantly lower than those at the unpaired codon sites are given a "+" sign

<sup>c</sup> datasets in which purifying selection detected within paired nucleotide sites was significantly stronger than that at unpaired nucleotide sites based on F and D statistics are given a "+" sign

<sup>d</sup> datasets in which a statistically significant association between paired sites and complementarily coevolving sites is detected are given a "+" sign